

# On stability of discretizations of the Helmholtz equation (extended version)

S. Esterhazy and J.M. Melenk

**Abstract** We review the stability properties of several discretizations of the Helmholtz equation at large wavenumbers. For a model problem in a polygon, a complete  $k$ -explicit stability (including  $k$ -explicit stability of the continuous problem) and convergence theory for high order finite element methods is developed. In particular, quasi-optimality is shown for a fixed number of degrees of freedom per wavelength if the mesh size  $h$  and the approximation order  $p$  are selected such that  $kh/p$  is sufficiently small and  $p = O(\log k)$ , and, additionally, appropriate mesh refinement is used near the vertices. We also review the stability properties of two classes of numerical schemes that use piecewise solutions of the homogeneous Helmholtz equation, namely, Least Squares methods and Discontinuous Galerkin (DG) methods. The latter includes the Ultra Weak Variational Formulation.

## 1 Introduction

A fundamental equation describing acoustic or electromagnetic phenomena is the time-dependent wave equation

$$\frac{\partial^2 w}{\partial t^2} - c^2 \Delta w = g,$$

given here for homogeneous, isotropic media whose propagation speed of waves is  $c$ . It arises in many applications, for example, radar/sonar detection, noise filtering, optical fiber design, medical imaging and seismic analysis. A commonly encountered setting is the time-harmonic case, in which the solution  $w$  (and correspondingly the right-hand side  $g$ ) is assumed to be of the form  $\operatorname{Re}(e^{-i\omega t} u(x))$  for

---

Vienna University of Technology, Institute for Analysis and Scientific Computing, Wiedner Hauptstrasse 8-10, A-1040 Vienna. e-mail: s.estershazy@tuwien.ac.at, melenk@tuwien.ac.at

a frequency  $\omega$ . Upon introducing the *wavenumber*  $k = \omega/c$  and the *wave length*  $\lambda := 2\pi/k$ , the resulting equation for the function  $u$ , which depends solely on the spatial variable  $x$ , is then the *Helmholtz equation*

$$-\Delta u - k^2 u = f. \quad (1)$$

In many high frequency situations of large  $k$  the solution  $u$  is highly oscillatory but has some multiscale character that can be captured, for example, by means of asymptotic analysis; a classical reference in this direction is [7].

In this article, we concentrate on numerical schemes for the Helmholtz equation at large wavenumbers  $k$ . Standard discretizations face several challenges, notably:

- (I) For large wavenumber  $k$ , the solution  $u$  is highly oscillatory. Its resolution, therefore, requires fine meshes, namely, at least  $N = k^d$  degrees of freedom, where  $d$  is the spatial dimension.
- (II) The standard  $H^1$ -conforming variational formulation is indefinite, and stability on the discrete level is therefore an issue. A manifestation of this problem is the so-called “pollution”, which expresses the observation that much more stringent conditions on the discretization have to be met than the minimal  $N = O(k^d)$  to achieve a given accuracy.

The second point, which will be the focus of the article, is best seen in the following, one-dimension example:

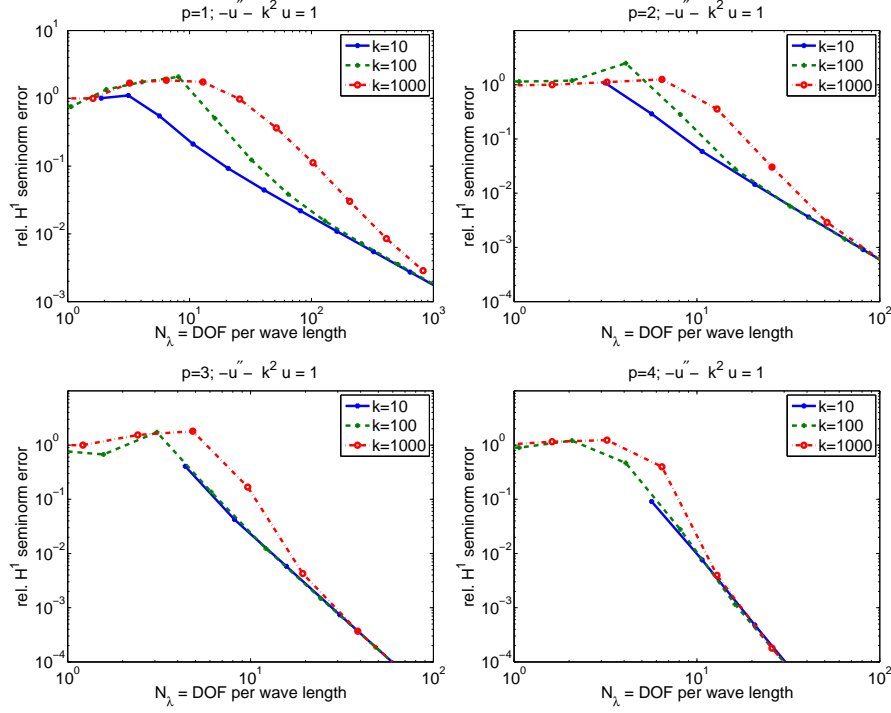
**Example 1.1.** For the boundary value problem

$$-u'' - k^2 u = 1 \quad \text{in } (0, 1), \quad u(0) = 0, \quad u'(1) - iku(1) = 0, \quad (2)$$

we consider the  $h$ -version finite element method (FEM) on uniform meshes with mesh size  $h$  for different approximation orders  $p \in \{1, 2, 3, 4\}$  and wavenumbers  $k \in \{1, 10, 100\}$ . Fig. 1 shows the relative error in the  $H^1(\Omega)$ -semi norm (i.e.,  $|u - u_N|_{H^1(\Omega)}/|u|_{H^1(\Omega)}$ , where  $u_N$  is the FEM approximation) versus the number of degrees of freedom per wavelength  $N_\lambda := N/\lambda = 2\pi N/k$  with  $N$  being the dimension of the finite element space employed. We observe several effects in Fig. 1: Firstly, since the solution  $u$  of (2) is smooth, higher order methods lead to higher accuracy for a given number of degrees of freedom per wavelength than lower order methods. Secondly, *asymptotically*, the FEM is quasioptimal with the finite element error  $|u - u_N|_{H^1(\Omega)}$  satisfying

$$|u - u_N|_{H^1(\Omega)} \approx C_p N_\lambda^{-p} |u|_{H^1(\Omega)} \quad (3)$$

for a constant  $C_p$  independent of  $k$ . Thirdly, the performance of the FEM as measured in “error vs. number of degrees of freedom per wavelength” does depend on  $k$ : As  $k$  increases, the preasymptotic range with reduced FEM performance becomes larger. Fourthly, higher order methods are less sensitive to  $k$  than lower order ones, i.e., for given  $k$ , high order methods enter the asymptotic regime in which (3) holds for smaller values of  $N_\lambda$  than lower order methods. ■



**Fig. 1** Performance of  $h$ -FEM for (2). Top:  $p = 1$ ,  $p = 2$ . Bottom:  $p = 3$ ,  $p = 4$  (see Example 1.1).

The behavior of the FEM in Example 1.1 has been analyzed in [47, 49], where error bounds of the form (see [47, Thm. 4.27])

$$|u - u_N|_{H^1(\Omega)} \leq C_p (1 + k^{p+1} h^p) h^p |u|_{H^{p+1}(\Omega)} \quad (4)$$

are established for a constant  $C_p$  depending only on the approximation order  $p$ . In this particular example, it is also easy to see that  $|u|_{H^{p+1}(\Omega)} / |u|_{H^1(\Omega)} \sim k^p$ , so that (4) can be recast in the form

$$|u - u_N|_{H^1(\Omega)} \leq C_p (1 + k^{p+1} h^p) (kh)^p |u|_{H^1(\Omega)} \sim (1 + kN_\lambda^{-p}) N_\lambda^{-p} |u|_{H^1(\Omega)}. \quad (5)$$

This estimate goes a long way to explain the above observations. The presence of the factor  $1 + kN_\lambda^{-p}$  explains the ‘‘pollution effect’’, i.e., the observation that for fixed  $N_\lambda$ , the (relative) error of the FEM as compared with the best approximation (which is essentially proportional to  $N_\lambda^{-p}$  in this example) increases with  $k$ . The estimate (5) also indicates that the asymptotic convergence behavior (3) is reached for  $N_\lambda = O(k^{1/p})$ . This confirms the observation made above that higher order methods are less prone to pollution than lower order methods. Although Example 1.1 is restricted to 1D, similar observations have been made in the literature also for

multi-d situations as early as [15]. We emphasize that for uniform meshes (as in Example 1.1) or, more generally, translation invariant meshes, complete and detailed dispersion analyses are available in an  $h$ -version setting, [2, 27, 47, 49], and in a  $p/hp$ -setting, [2–4], that give strong mathematical evidence for the superior ability of high order methods to cope with the pollution effect.

The present paper, which discusses and generalizes the work [61, 62], proves that also on unstructured meshes, high order methods are less prone to pollution. More precisely, for a large class of Helmholtz problems, stability and quasi-optimality is given under the scale resolution condition

$$\frac{kh}{p} \leq c_1 \quad \text{and} \quad p \geq c_2 \log k, \quad (6)$$

where  $c_1$  is sufficiently small and  $c_2$  sufficiently large. For piecewise smooth geometries (e.g., polygons), additionally appropriate mesh refinement near the singularities is required.

We close our discussion of Example 1.1 by remarking that its analysis and, in fact, the analysis of significant parts of this article rests on  $H^1$ -like norms. Largely, this choice is motivated by the numerical scheme, namely, an  $H^1$ -conforming FEM.

### 1.1 Non-standard FEM

The limitations of the classical FEM mentioned above in (I) and (II) have sparked a significant amount of research in the past decades to overcome or at least mitigate them. This research focuses on two techniques that are often considered in tandem: firstly, the underlying approximation by classical piecewise polynomials is replaced with special, problem-adapted functions such as systems of plane waves; secondly, the numerical scheme is based on a variational formulation different from the classical  $H^1$ -conforming Galerkin approach. Before discussing these ideas in more detail, we point the reader to the interesting work [12], which shows for a model situation on regular, infinite grids in 2D that no 9-point stencil (i.e., a numerical method based on connecting the value at a node with those of the 8 nearest neighbors) generates a completely pollution-free method; the 1D situation is special and discussed briefly in Section 7.

Work that is based on a new or modified variational formulation but rests on the approximation properties of piecewise polynomials includes the Galerkin Least Squares Method [39, 40], the methods of [9], and Discontinuous Galerkin Methods ([33–35] and references there). Several methods have been proposed that are based on the approximation properties of special, problem-adapted systems of functions such as systems of plane waves. In an  $H^1$ -conforming Galerkin setting, this idea has been pursued in the Partition of Unity Method/Generalized FEM by several authors, e.g., [5, 45, 50, 51, 56, 60, 68, 69, 81]. A variety of other methods that are based on problem-adapted ansatz functions leave the  $H^1$ -conforming Galerkin setting and en-

force the jump across element boundaries in a weak sense. This can be done by least squares techniques ([14, 53, 65, 70, 80] and references there), by Lagrange multiplier techniques as in the Discontinuous Enrichment Method [31, 32, 82] or by Discontinuous Galerkin (DG) type methods, [19–21, 36, 42, 43, 46, 55, 63, 64]; in these last references, we have included the work on the Ultra Weak Variational Formulation (UWVF) since it can be understood as a special DG method as discussed in [19, 36].

## 1.2 Scope of the article

The present article focuses on the stability properties of numerical methods for Helmholtz problems and exemplarily discusses three different approaches in more detail for their differences in techniques. The first approach, studied in Section 4, is that of the classical  $H^1$ -based Galerkin method for Helmholtz problems. The setting is that of a Gårding inequality so that stability of a numerical method can be inferred from the stability of the continuous problem by perturbation arguments. This motivates us to study for problem (9), which will serve as our model Helmholtz problem in this article, the stability properties of the continuous problem in Section 2. In order to make the perturbation argument explicit in the wavenumber  $k$ , a detailed,  $k$ -explicit regularity analysis for Helmholtz problems is necessary. This is worked out in Section 3 for our model problem (9) posed on polygonal domains. These results generalize a similar regularity theory for convex polygons or domains with analytic boundary of [61, 62]. Structurally similar results have been obtained in connection with boundary integral formulations in [54, 59].

We discuss in Sections 6.2 and 6.3 somewhat briefly a second and a third approach to stability of numerical schemes. In contrast to the setting discussed above, where stability is only ensured asymptotically for sufficiently fine discretizations, these methods are stable by construction and can even feature quasioptimality in appropriate residual norms. We point out, however, that relating this residual norm to a more standard norm such as the  $L^2$ -norm for the error is a non-trivial task. Our presentation for these methods will follow the works [19, 36, 43, 65].

Many aspects of discretizations for Helmholtz problems are not addressed in this article. For recent developments in boundary element techniques for this problem class, we refer to the survey article [22]. The model problem (9) discussed here involves the rather simple boundary condition (9b), which can be understood as an approximation to a Dirichlet-to-Neumann operator that provides a coupling to a homogeneous Helmholtz equation in an exterior domain together with appropriate radiation conditions at infinity. A variety of techniques for such problems are discussed in [37]. Further methods include FEM-BEM coupling, the PML due to Bérenger (see [17, 24] and references therein), infinite elements [26], and methods based on the pole condition, [44]. Another topic not addressed here is the solution of the arising linear system; we refer the reader to [28, 30] for a discussion of the state of the art. Further works with survey character includes [29, 47, 48, 83].

### 1.3 Some notation

We employ standard notation for Sobolev spaces, [1, 18, 67, 77]. For domains  $\omega$  and  $k \neq 0$  we denote

$$\|u\|_{1,k,\omega}^2 := k^2 \|u\|_{L^2(\omega)}^2 + \|\nabla u\|_{L^2(\omega)}^2. \quad (7)$$

This norm is equivalent to the standard  $H^1$ -norm. The presence of the weight  $k$  in the  $L^2$ -part leads to a balance between the  $H^1$ -seminorm and the  $L^2$ -norm for functions with the expected oscillatory behavior such as plane waves  $e^{i\mathbf{k}\mathbf{d}\cdot\mathbf{x}}$  (with  $\mathbf{d}$  being a unit vector). Additionally, the bilinear form  $B$  considered below is bounded uniformly in  $k$  with respect to this ( $k$ -dependent) norm.

Throughout this work, a standing assumption will be

$$|k| \geq k_0 > 0; \quad (8)$$

our frequently used phrase “independent of  $k$ ” will still implicitly assume (8). We denote by  $C$  a generic constant. If not stated otherwise,  $C$  will be independent of the wavenumber  $k$  but may depend on  $k_0$ . For smooth functions  $u$  defined on a  $d$ -dimensional manifold, we employ the notation  $|\nabla^n u(x)|^2 := \sum_{\alpha \in \mathbb{N}_0^d: |\alpha|=n} \frac{|\alpha|!}{\alpha!} |D^\alpha u(x)|^2$ .

### 1.4 A model problem

In order to fix ideas, we will use the following, specific model problem: For a bounded Lipschitz domain  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , we study for  $k \in \mathbb{R}$ ,  $|k| \geq k_0$ , the boundary value problem

$$-\Delta u - k^2 u = f \text{ in } \Omega, \quad (9a)$$

$$\partial_n u + iku = g \text{ on } \partial\Omega. \quad (9b)$$

Henceforth, to simplify the notation, we assume  $k \geq k_0 > 0$  but point out that the choice of the sign of  $k$  is not essential. The weak formulation for (9) is:

$$\text{Find } u \in H^1(\Omega) \text{ s.t. } B(u, v) = l(v) \quad \forall v \in H^1(\Omega), \quad (10)$$

where, for  $f \in L^2(\Omega)$  and  $g \in L^2(\partial\Omega)$ ,  $B$  and  $l$  are given by

$$B(u, v) := \int_{\Omega} (\nabla u \cdot \nabla \bar{v} - k^2 u \bar{v}) + i k \int_{\partial\Omega} u \bar{v}, \quad l(v) := (f, v)_{L^2(\Omega)} + (g, v)_{L^2(\partial\Omega)}. \quad (11)$$

As usual, if  $f \in (H^1(\Omega))'$  and  $g \in H^{-1/2}(\partial\Omega)$ , then the  $L^2$ -inner products  $(\cdot, \cdot)_{L^2(\Omega)}$  and  $(\cdot, \cdot)_{L^2(\partial\Omega)}$  are understood as duality pairings. The multiplicative trace inequality proves continuity of  $B$ ; in fact, there exists  $C_B > 0$  independent of  $k$  such that (see, e.g., [61, Cor. 3.2] for details)

$$|B(u, v)| \leq C_B \|u\|_{1,k,\Omega} \|v\|_{1,k,\Omega} \quad \forall u, v \in H^1(\Omega). \quad (12)$$

## 2 Stability of the continuous problem

Helmholtz problems can often be cast in the form “coercive + compact perturbation” where the compact perturbation is  $k$ -dependent. In other words, a Gårding inequality is satisfied. For example, the sesquilinear form  $B$  of (11) is of this form since

$$\operatorname{Re} B(u, u) + 2k^2(u, u)_{L^2(\Omega)} = \|u\|_{1,k,\Omega}^2 \quad (13)$$

and the embedding  $H^1(\Omega) \subset L^2(\Omega)$  is compact by Rellich’s theorem. Classical Fredholm theory (the “Fredholm alternative”) then yields unique solvability of (10) for all  $f \in (H^1(\Omega))'$  and  $g \in H^{-1/2}(\partial\Omega)$ , if one can show uniqueness. Uniqueness in turn is often obtained by exploiting analyticity of the solutions of homogeneous Helmholtz equation, or, more generally, the unique continuation principle for elliptic problems, (see, e.g., [52, Chap. 4.3]):

**Example 2.1 (Uniqueness for (9)).** Let  $f = 0$  and  $g = 0$  in (9). Then, any solution  $u \in H^1(\Omega)$  of (9) satisfies  $u|_{\partial\Omega} = 0$  since  $0 = \operatorname{Im} B(u, u) = k \|u\|_{L^2(\partial\Omega)}^2$  (see Lemma 2.2). Hence, the trivial extension  $\tilde{u}$  to  $\mathbb{R}^2$  satisfies  $\tilde{u} \in H^1(\mathbb{R}^2)$ . The observations  $B(u, v) = 0$  for all  $v \in H^1(\Omega)$  and  $u|_{\partial\Omega} = 0$  show

$$\int_{\mathbb{R}^2} \nabla \tilde{u} \cdot \nabla \bar{v} - k^2 \tilde{u} \bar{v} = 0 \quad \forall v \in C_0^\infty(\mathbb{R}^2).$$

Hence,  $\tilde{u}$  is a solution of the homogeneous Helmholtz equation and  $\tilde{u}$  vanishes on  $\mathbb{R}^2 \setminus \overline{\Omega}$ . Analyticity of  $\tilde{u}$  (or, more generally, the unique continuation principle presented in [52, Chap. 4.3]) then implies that  $\tilde{u} \equiv 0$ , which in turn yields  $u \equiv 0$ . ■

The arguments based on the Fredholm alternative do not give any indication of how the solution operator depends on the wavenumber  $k$ . Yet, it is clearly of interest to know how  $k$  enters bounds for the solution operator. It turns out that both the geometry and the type of boundary conditions strongly affect these bounds. For example, for an exterior Dirichlet problem, [16] exhibits a geometry and a sequence of wavenumber  $(k_n)_{n \in \mathbb{N}}$  tending to infinity such that the norm of the solution operator for these wavenumbers is bounded from below by an exponentially growing term  $Ce^{bk_n}$  for some  $C, b > 0$ . These geometries feature so-called “trapping” or near-trapping and are not convex. For convex or at least star-shaped geometries, the  $k$ -dependence is much better behaved. An important ingredient of the analysis on such geometries are special test functions in the variational formulation. For example, assuming in the the model problem (10) that  $\Omega$  is star-shaped with respect to the origin (and has a smooth boundary), one may take as the test function  $v(x) = x \cdot \nabla u(x)$ , where  $u$  is the exact solution. An integration by parts (more generally, the so-called “Rellich identities” [67, p. 261] or an identity due to Pohožaev, [71]) then leads to the following estimate for the model problem (10):

$$\|u\|_{1,k,\Omega} \leq C \left[ \|f\|_{L^2(\Omega)} + \|g\|_{L^2(\partial\Omega)} \right]; \quad (14)$$

this was shown in [56, Prop. 8.1.4] (for  $d = 2$ ) and subsequently by [25] for  $d = 3$ . Uniform in  $k$  bounds were established in [41] for star-shaped domains and certain boundary conditions of mixed type by related techniques. The same test function was also crucial for a boundary integral setting in [23]. A refined version of this test function that goes back to Morawetz and Ludwig, [66] was used recently in a boundary integral equations context (still for star-shaped domains), [78].

While (14) does not make minimal assumptions on the regularity of  $f$  and  $g$ , the estimate (14) can be used to show that (for star-shaped domains) the sesquilinear form  $B$  of (10) satisfies an inf-sup condition with inf-sup constant  $\gamma = O(k^{-1})$ —this can be shown using the arguments presented in the proof Theorem 2.5.

An important ingredient of the regularity and stability theory will be the concept of *polynomial well-posedness* by which we mean polynomial-in- $k$ -bounds for the norm of the solution operator. The model problem (9) on star-shaped domains with the *a priori* bound (14) is an example. The following Section 2.1 shows polynomial well-posedness for the model problem (9) on general Lipschitz domains (Thm. 2.4). It is thus not the geometry but the type of boundary conditions in our model problem (9), namely, Robin boundary conditions that makes it polynomially well-posed. In contrast, the Dirichlet boundary conditions in conjunction with the lack of star-shapedness in the examples given in [16] make these problem not polynomially well-posed.

## 2.1 Polynomial well-posedness for the model problem (9)

**Lemma 2.2.** *Let  $\Omega \subset \mathbb{R}^d$  be a bounded Lipschitz domain. Let  $u \in H^1(\Omega)$  be a weak solution of (9) with  $f = 0$  and  $g \in L^2(\partial\Omega)$ . Then  $\|u\|_{L^2(\partial\Omega)} \leq k^{-1} \|g\|_{L^2(\partial\Omega)}$ .*

*Proof.* Selecting  $v = u$  in the weak formulation (10) and considering the imaginary part yields

$$k \|u\|_{L^2(\partial\Omega)}^2 = \operatorname{Im} \int_{\partial\Omega} g \bar{u} \leq \|g\|_{L^2(\partial\Omega)} \|u\|_{L^2(\partial\Omega)}.$$

This concludes the argument.  $\square$

Next we use results on layer potentials for the Helmholtz equation from [59] to prove the following lemma:

**Lemma 2.3.** *Let  $\Omega \subset \mathbb{R}^d$  be a bounded Lipschitz domain,  $u \in H^1(\Omega)$  solve (9) with  $f = 0$ . Assume  $u|_{\partial\Omega} \in L^2(\partial\Omega)$  and  $\partial_n u \in L^2(\partial\Omega)$ . Then there exists  $C > 0$  independent of  $k$  and  $u$  such that*

$$\begin{aligned} \|u\|_{L^2(\Omega)} &\leq Ck \left( \|u\|_{L^2(\partial\Omega)} + \|\partial_n u\|_{H^{-1}(\partial\Omega)} \right), \\ \|u\|_{1,k,\Omega} &\leq C \left[ k^2 \|u\|_{L^2(\partial\Omega)} + k^2 \|\partial_n u\|_{H^{-1}(\partial\Omega)} + k^{-2} \|\partial_n u\|_{L^2(\partial\Omega)} \right]. \end{aligned}$$



*Proof.* With the single layer and double layer potentials  $\tilde{V}_k$  and  $\tilde{K}_k$  we have the representation formula  $u = \tilde{V}_k \partial_n u - \tilde{K}_k u$ . From [59, Lemmata 2.1, 2.2, Theorems 4.1, 4.2] we obtain

$$\|\tilde{V}_k \partial_n u\|_{L^2(\Omega)} \leq Ck \|\partial_n u\|_{H^{-1}(\partial\Omega)}, \quad \|\tilde{K}_k u\|_{L^2(\Omega)} \leq Ck \|u\|_{L^2(\partial\Omega)}.$$

Thus,

$$\|u\|_{L^2(\Omega)} \leq Ck \left( \|u\|_{L^2(\partial\Omega)} + \|\partial_n u\|_{H^{-1}(\partial\Omega)} \right).$$

Next, using  $v = u$  in the weak formulation (10) yields

$$\|\nabla u\|_{L^2(\Omega)}^2 \leq C \left[ k^2 \|u\|_{L^2(\Omega)}^2 + \|\partial_n u\|_{L^2(\partial\Omega)} \|u\|_{L^2(\partial\Omega)} \right]$$

and therefore

$$\|\nabla u\|_{L^2(\Omega)}^2 + k^2 \|u\|_{L^2(\Omega)}^2 \leq C \left[ k^4 \|u\|_{L^2(\partial\Omega)}^2 + k^4 \|\partial_n u\|_{H^{-1}(\partial\Omega)}^2 + k^{-4} \|\partial_n u\|_{L^2(\partial\Omega)}^2 \right],$$

which concludes the proof.  $\square$

**Theorem 2.4.** *Let  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$  be a bounded Lipschitz domain. Then there exists  $C > 0$  (independent of  $k$ ) such that for  $f \in L^2(\Omega)$  and  $g \in L^2(\partial\Omega)$  the solution  $u \in H^1(\Omega)$  of (9) satisfies*

$$\|u\|_{1,k,\Omega} \leq C \left[ k^2 \|g\|_{L^2(\partial\Omega)} + k^{5/2} \|f\|_{L^2(\Omega)} \right].$$

*Proof.* We first transform the problem to one with homogeneous right-hand side  $f$  in the standard way. A particular solution of the equation (9a) is given by the Newton potential  $u_0 := G_k \star f$ ; here,  $G_k$  is a Green's function for the Helmholtz equation and we tacitly extend  $f$  by zero outside  $\Omega$ . Then  $u_0 \in H_{loc}^2(\mathbb{R}^d)$  and by the analysis of the Newton potential given in [61, Lemma 3.5] we have

$$k^{-1} \|u_0\|_{H^2(\Omega)} + \|u_0\|_{H^1(\Omega)} + k \|u_0\|_{L^2(\Omega)} \leq C \|f\|_{L^2(\Omega)}. \quad (15)$$

The difference  $\tilde{u} := u - u_0$  then satisfies

$$-\Delta \tilde{u} - k^2 \tilde{u} = 0 \quad \text{in } \Omega, \quad (16a)$$

$$\partial_n \tilde{u} + ik \tilde{u} = g - (\partial_n u_0 + ik u_0) =: \tilde{g}. \quad (16b)$$

We have with the multiplicative trace inequality

$$\begin{aligned} \|\tilde{g}\|_{L^2(\partial\Omega)} &\leq C \left[ \|g\|_{L^2(\partial\Omega)} + \|u_0\|_{H^2(\Omega)}^{1/2} \|u_0\|_{H^1(\Omega)}^{1/2} + k \|u_0\|_{H^1(\Omega)}^{1/2} \|u_0\|_{L^2(\Omega)}^{1/2} \right] \\ &\leq C \left[ \|g\|_{L^2(\partial\Omega)} + k^{1/2} \|f\|_{L^2(\Omega)} \right]. \end{aligned} \quad (17)$$

To get bounds on  $\tilde{u}$ , we employ Lemma 2.2 and (17) to conclude

$$\|\tilde{u}\|_{L^2(\partial\Omega)} \leq Ck^{-1}\|\tilde{g}\|_{L^2(\partial\Omega)} \leq C \left[ k^{-1}\|g\|_{L^2(\partial\Omega)} + k^{-1/2}\|f\|_{L^2(\Omega)} \right], \quad (18)$$

$$\|\partial_n \tilde{u}\|_{L^2(\partial\Omega)} \leq C \left[ \|\tilde{g}\|_{L^2(\partial\Omega)} + k\|\tilde{u}\|_{L^2(\partial\Omega)} \right] \leq C \left[ \|g\|_{L^2(\partial\Omega)} + k^{1/2}\|f\|_{L^2(\Omega)} \right]. \quad (19)$$

Lemma 2.3 and the generous estimate  $\|\partial_n \tilde{u}\|_{H^{-1}(\partial\Omega)} \leq C\|\partial_n \tilde{u}\|_{L^2(\partial\Omega)}$  produce

$$\|\tilde{u}\|_{H^1(\Omega)} + k\|\tilde{u}\|_{L^2(\Omega)} \leq C \left[ k^2\|g\|_{L^2(\partial\Omega)} + k^{5/2}\|f\|_{L^2(\Omega)} \right]. \quad (20)$$

Combining (15), (20) finishes the argument.  $\square$

The *a priori* estimate of Theorem 2.4 does not make minimal assumptions on the regularity of  $f$  and  $g$ . However, it can be used to obtain estimates on the inf-sup and hence *a priori* bounds for  $f \in (H^1(\Omega))'$  and  $g \in H^{-1/2}(\partial\Omega)$  as we now show:

**Theorem 2.5.** *Let  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$  be a bounded Lipschitz domain. Then there exists  $C > 0$  (independent of  $k$ ) such that the sesquilinear form  $B$  of (11) satisfies*

$$\inf_{0 \neq u \in H^1(\Omega)} \sup_{0 \neq v \in H^1(\Omega)} \frac{\operatorname{Re} B(u, v)}{\|u\|_{1,k,\Omega} \|v\|_{1,k,\Omega}} \geq Ck^{-7/2}. \quad (21)$$

Furthermore, for every  $f \in (H^1(\Omega))'$  and  $g \in H^{-1/2}(\partial\Omega)$  the problem (10) is uniquely solvable, and its solution  $u \in H^1(\Omega)$  satisfies the *a priori* bound

$$\|u\|_{1,k,\Omega} \leq Ck^{7/2} \left[ \|f\|_{(H^1(\Omega))'} + \|g\|_{H^{-1/2}(\partial\Omega)} \right]. \quad (22)$$

If  $\Omega$  is convex or if  $\Omega$  is star-shaped and has a smooth boundary, then the following, sharper estimate holds:

$$\inf_{0 \neq u \in H^1(\Omega)} \sup_{0 \neq v \in H^1(\Omega)} \frac{\operatorname{Re} B(u, v)}{\|u\|_{1,k,\Omega} \|v\|_{1,k,\Omega}} \geq Ck^{-1}. \quad (23)$$

*Proof.* The proof relies on standard arguments for sesquilinear forms satisfying a Gårding inequality. For simplicity of notation, we write  $\|\cdot\|_{1,k}$  for  $\|\cdot\|_{1,k,\Omega}$ .

Given  $u \in H^1(\Omega)$  we define  $z \in H^1(\Omega)$  as the solution of

$$2k^2(\cdot, u)_{L^2(\Omega)} = B(\cdot, z).$$

Theorem 2.4 implies  $\|z\|_{1,k} \leq Ck^{9/2}\|u\|_{L^2(\Omega)}$ , and  $v = u + z$  satisfies

$$\operatorname{Re} B(u, v) = \operatorname{Re} B(u, u) + \operatorname{Re} B(u, z) = \|u\|_{1,k}^2 - 2k^2\|u\|_{L^2(\Omega)}^2 + \operatorname{Re} B(u, z) = \|u\|_{1,k}^2.$$

Thus,

$$\operatorname{Re} B(u, v) = \|u\|_{1,k}^2,$$

$$\|v\|_{1,k} = \|u + z\|_{1,k} \leq \|u\|_{1,k} + \|z\|_{1,k} \leq \|u\|_{1,k} + Ck^{9/2}\|u\|_{L^2(\Omega)} \leq Ck^{7/2}\|u\|_{1,k}.$$

Therefore,

$$\operatorname{Re} B(u, v) = \|u\|_{1,k}^2 \geq \|u\|_{1,k} Ck^{-7/2} \|v\|_{1,k},$$

which concludes the proof of (21). Example 2.1 provides unique solvability for (9) so that (21) gives the *a priori* estimate (22). Finally, (23) is shown by the same arguments using (14).  $\square$

### 3 $k$ -explicit regularity theory

#### 3.1 Regularity by decomposition

Since the Sobolev regularity of elliptic problems is determined by the leading order terms of the differential equation and the boundary conditions, the Sobolev regularity properties of our model problem (9) are the same as those for the Laplacian. However, regularity results that are explicit in the wavenumber  $k$  are clearly of interest; for example, we will use them in Section 4.2 below to quantify how fine the discretization has to be (relative to  $k$ ) so that the FEM is stable and quasi-optimal.

The  $k$ -explicit regularity theory developed in [61, 62] (and, similarly, for integral equations in [54, 59]) takes the form of an additive splitting of the solution into a part with finite regularity but  $k$ -independent bounds and a part that is analytic and for which  $k$ -explicit bounds for all derivatives are available. Below, we will present a similar regularity theory for the model problem (9) for polygonal  $\Omega \subset \mathbb{R}^2$ , thereby extending the results of [62], which restricted its analysis of polygons to the convex case. In order to motivate the ensuing developments, we quote from [61] a result that shows in a simple setting the key features of our  $k$ -explicit “regularity by decomposition”:

**Lemma 3.1 ([61, Lemma 3.5]).** *Let  $B_R(0) \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$  be the ball of radius  $R$  centered at the origin. Then, there exist  $C, \gamma > 0$  such that for all  $k$  (with  $k \geq k_0$ ) the following is true: For all  $f \in L^2(\mathbb{R}^d)$  with  $\operatorname{supp} f \subset B_R(0)$  the solution  $u$  of*

$$-\Delta u - k^2 u = f \quad \text{in } \mathbb{R}^d,$$

*subject to the Sommerfeld radiation condition*

$$\lim_{|x| \rightarrow \infty} |x|^{\frac{d-1}{2}} \left( \frac{\partial u}{\partial |x|} - iku \right) = 0 \quad \text{for } |x| \rightarrow \infty,$$

*has the following regularity properties:*

- (i)  $u|_{B_{2R}(0)} \in H^2(B_{2R}(0))$  and  $\|u\|_{H^2(B_{2R}(0))} \leq Ck\|f\|_{L^2(B_R(0))}$ .
- (ii)  $u|_{B_{2R}(0)}$  can be decomposed as  $u = u_{H^2} + u_{\mathcal{A}}$  for a  $u_{H^2} \in H^2(B_{2R})$  and an analytic  $u_{\mathcal{A}}$  together with the bounds

$$\begin{aligned}
k\|u_{H^2}\|_{1,k,B_{2R}(0)} + \|u_{H^2}\|_{H^2(B_{2R}(0))} &\leq C\|f\|_{L^2(B_R(0))}, \\
\|\nabla^n u_{\mathcal{A}}\|_{L^2(B_{2R}(0))} &\leq C\gamma^n \max\{n, k\}^{n-1} \|f\|_{L^2(B_R(0))} \quad \forall n \in \mathbb{N}_0.
\end{aligned}$$

A few comments concerning Lemma 3.1 are in order. For general  $f \in L^2(B_R(0))$ , one cannot expect better regularity than  $H^2$ -regularity for the solution  $u$  and, indeed, both regularity results (i) and (ii) assert this. The estimate (i) is sharp in its dependence on  $k$  as the following simple example shows: For the fundamental solution  $G_k$  (with singularity at the origin) and a cut-off function  $\chi \in C_0^\infty(\mathbb{R}^d)$  with  $\text{supp } \chi \subset B_{2R}(0)$  and  $\chi \equiv 1$  on  $B_R(0)$ , the functions  $u := (1 - \chi)G_k$  and  $f := -\Delta u - k^2 u$  satisfy  $\|u\|_{H^2(B_{2R}(0))} = O(k^2)$  and  $\|f\|_{L^2(B_R(0))} = O(k)$ . Compared to (i), the regularity assertion (ii) is finer in that its  $H^2$ -part  $u_{H^2}$  has a better  $k$ -dependence. The  $k$ -dependence of the analytic part  $u_{\mathcal{A}}$  is not improved (indeed,  $\|u_{\mathcal{A}}\|_{H^2(B_{2R}(0))} \leq Ck\|f\|_{L^2(B_R(0))}$ ), but the analyticity of  $u_{\mathcal{A}}$  is a feature that higher order methods can exploit.

The decomposition in (ii) of Lemma 3.1 is obtained by a decomposition of the datum  $f$  using low pass and high pass filters, i.e.,  $f = L_{\eta k}f + H_{\eta k}f$ , where the low pass filter  $L_{\eta k}$  cuts off frequencies beyond  $\eta k$  (here,  $\eta > 1$ ) and  $H_{\eta k}$  eliminates the frequencies small than  $\eta k$ . Similar frequency filters will be important tools in our analysis below as well (see Sec. 3.3.1). The regularity properties stated in (ii) then follow from this decomposition and the explicit solution formula  $u = G_k \star f$  (see [61, Lemma 3.5] for details).

Lemma 3.1 serves as a prototype for “regularity theory by decomposition”. Similar decompositions have been developed recently for several Helmholtz problems in [62] and [54, 59]. Although they vary in their details, these decomposition are structurally similar in that they have the form of an additive splitting into a part with finite regularity with  $k$ -independent bounds and an analytic part with  $k$ -dependent bounds. The basic ingredients of these decomposition results are (a) (piecewise) analyticity of the geometry (or, more generally, the data of the problem) and (b) *a priori* bounds for solution operator. The latter appear only in the estimate for the analytic part of the decomposition, and the most interesting case is that of polynomially well-posed problems. We illustrate the construction of the decomposition for the model problem (9) in polygonal domains  $\Omega \subset \mathbb{R}^2$ . This result is an extension to general polygons of the results [62], which restricted its attention to the case of convex polygons. We emphasize that the choice of the boundary conditions (9b) is not essential for the form of the decomposition and other boundary conditions could be treated using similar techniques.

### 3.2 Setting and main result

Let  $\Omega \subset \mathbb{R}^2$  be a bounded, polygonal Lipschitz domain with vertices  $A_j$ ,  $j = 1, \dots, J$ , and interior angles  $\omega_j$ ,  $j = 1, \dots, J$ . We will require the countably normed spaces introduced in [8, 57]. These space are designed to capture important features of solu-

tions of elliptic partial differential equations posed on polygons, namely, analyticity of the solution and the singular behavior at the vertices. Their characterization in terms of these countably normed spaces also permits proving exponential convergence of piecewise polynomial approximation on appropriately graded meshes.

These countably normed spaces are defined with the aid of weight functions  $\Phi_{p,\vec{\beta},k}$  that we now define. For  $\beta \in [0, 1)$ ,  $n \in \mathbb{N}_0$ ,  $k > 0$ , and  $\vec{\beta} \in [0, 1)^J$ , we set

$$\Phi_{n,\beta,k}(x) = \min \left\{ 1, \frac{|x|}{\min \left\{ 1, \frac{|n|+1}{k+1} \right\}} \right\}^{n+\beta},$$

$$\Phi_{n,\vec{\beta},k}(x) = \prod_{j=1}^J \Phi_{n,\beta_j,k}(x - A_j). \quad (24)$$

Finally, we denote by  $H_{pw}^{1/2}(\partial\Omega)$  the space of functions whose restrictions of the edges of  $\partial\Omega$  are in  $H^{1/2}$ .

We furthermore introduce the constant  $C_{sol}(k)$  as a suitable norm of the solution operator for (9). That is,  $C_{sol}(k)$  is such that for all  $f \in L^2(\Omega)$ ,  $g \in L^2(\partial\Omega)$  and corresponding solution  $u$  of (9) there holds

$$\|u\|_{1,k,\Omega} \leq C_{sol}(k) \left[ \|f\|_{L^2(\Omega)} + \|g\|_{L^2(\partial\Omega)} \right]. \quad (25)$$

We recall that Theorem 2.4 gives  $C_{sol}(k) = O(k^{5/2})$  for general polygons and  $C_{sol}(k) = O(1)$  by [56, Prop. 8.1.4] for convex polygons. Our motivation for using the notation  $C_{sol}(k)$  is emphasize in the following theorem how *a priori* estimates for Helmholtz problems affect the decomposition result:

**Theorem 3.2.** *Let  $\Omega \subset \mathbb{R}^2$  be a polygon with vertices  $A_j$ ,  $j = 1, \dots, J$ . Then there exist constants  $C$ ,  $\gamma > 0$ ,  $\vec{\beta} \in [0, 1)^J$  independent of  $k \geq k_0$  such that for every  $f \in L^2(\Omega)$  and  $g \in H_{pw}^{1/2}(\partial\Omega)$  the solution  $u$  of (9) can be written as  $u = u_{H^2} + u_{\mathcal{A}}$  with*

$$\begin{aligned} k\|u_{H^2}\|_{1,k,\Omega} + \|u_{H^2}\|_{H^2(\Omega)} &\leq CC_{f,g} \\ \|u_{\mathcal{A}}\|_{H^1(\Omega)} &\leq (C_{sol}(k) + 1)C_{f,g} \\ k\|u_{\mathcal{A}}\|_{L^2(\Omega)} &\leq (C_{sol}(k) + k)C_{f,g} \\ \|\Phi_{n,\vec{\beta},k} \nabla^{n+2} u_{\mathcal{A}}\|_{L^2(\Omega)} &\leq C(C_{sol}(k) + 1)k^{-1} \max\{n, k\}^{n+2} \gamma^n C_{f,g} \quad \forall n \in \mathbb{N}_0 \end{aligned}$$

with  $C_{f,g} := \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)}$  and  $C_{sol}(k)$  introduced in (25).

*Proof.* The proof is relegated to Section 3.4. We mention that the  $k$ -dependence of our bounds on  $\|u_{\mathcal{A}}\|_{L^2(\Omega)}$  is likely to be suboptimal due to our method of proof.  $\square$

Theorem 3.2 may be viewed as the analog of Lemma 3.1, (ii); we conclude this section with the analog of Lemma 3.1, (i):

**Corollary 3.3.** *Assume the hypotheses of Theorem 3.2. Then there exist constants  $C > 0$ ,  $\vec{\beta} \in [0, 1]^J$  independent of  $k$  such that for all  $f \in L^2(\Omega)$ ,  $g \in H_{pw}^{1/2}(\partial\Omega)$  the solution  $u$  of (9) satisfies  $\|u\|_{1,k,\Omega} \leq CC_{sol}(k) \left[ \|f\|_{L^2(\Omega)} + \|g\|_{L^2(\partial\Omega)} \right]$  as well as*

$$\|\Phi_{0,\vec{\beta},k} \nabla^2 u\|_{L^2(\Omega)} \leq Ck(C_{sol}(k) + 1) \left[ \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)} \right].$$

*Proof.* The estimate for  $\|u\|_{1,k,\Omega}$  expresses (25). The estimate for the second derivatives of  $u$  follows from Theorem 3.2 since  $u = u_{H^2} + u_{\mathcal{A}}$ .  $\square$

### 3.3 Auxiliary results

Just as in the proof of Lemma 3.1, an important ingredient of the proof of Theorem 3.2 are high and low pass filters. The underlying reason is that the Helmholtz operator  $-\Delta - k^2$  acts very differently on low frequency and high frequency functions. Here, the dividing line between high and low frequencies is at  $O(k)$ . For this reason, appropriate high and low pass filters are defined and analyzed in Section 3.3.1. Furthermore, when applied to high frequency functions the Helmholtz operator behaves similarly to the Laplacian  $-\Delta$  or the modified Helmholtz operator  $-\Delta + k^2$ . This latter operator, being positive definite, is easier to analyze and yet provides insight into the behavior of the Helmholtz operator restricted to high frequency functions. The modified Helmholtz operator will therefore be a tool for the proof of Theorem 3.2 and is thus analyzed in Section 3.3.3.

#### 3.3.1 High and low pass filters, auxiliary results

For the polygonal domain  $\Omega \subset \mathbb{R}^2$  we introduce for  $\eta > 1$  the following two low and high pass filters in terms of the Fourier transform  $\mathcal{F}$ :

1. The low and high pass filters  $L_{\Omega,\eta}f : L^2(\Omega) \rightarrow L^2(\Omega)$  and  $H_{\Omega,\eta} : L^2(\Omega) \rightarrow L^2(\Omega)$  are defined by

$$L_{\Omega,\eta}f = (\mathcal{F}^{-1} \chi_{B_{\eta k}(0)} \mathcal{F}(E_{\Omega}f))|_{\Omega}, \quad H_{\Omega,\eta}f = (\mathcal{F}^{-1} \chi_{\mathbb{R}^2 \setminus B_{\eta k}(0)} \mathcal{F}(E_{\Omega}f))|_{\Omega};$$

here,  $B_{\eta k}(0)$  is the ball of radius  $\eta k$  with center 0, the characteristic function of a set  $A$  is  $\chi_A$ , and  $E_{\Omega}$  denotes the Stein extension operator of [79, Chap. VI].

2. Analogously, we define  $L_{\partial\Omega,\eta}f : L^2(\partial\Omega) \rightarrow L^2(\partial\Omega)$  and  $H_{\partial\Omega,\eta} : L^2(\partial\Omega) \rightarrow L^2(\partial\Omega)$  in an *edgewise* fashion. Specifically, identifying an edge  $e$  of  $\Omega$  with an interval and letting  $E_e$  be the Stein extension operator for the interval  $e \subset \mathbb{R}$  to the real line  $\mathbb{R}$ , we can define with the univariate Fourier transformation  $\mathcal{F}$  the operators  $L_{e,\eta}$  and  $H_{e,\eta}$  by

$$L_{e,\eta}g = (\mathcal{F}^{-1}\chi_{B_{\eta k}(0)}\mathcal{F}(E_e g))|_e, \quad H_{e,\eta}g = (\mathcal{F}^{-1}\chi_{\mathbb{R}\setminus B_{\eta k}(0)}\mathcal{F}(E_e f))|_e;$$

the operators  $L_{\partial\Omega,\eta}$  and  $H_{\partial\Omega,\eta}$  are then defined edgewise by  $(L_{\partial\Omega,\eta}g)|_e = L_{e,\eta}g$  and  $(H_{\partial\Omega,\eta}g)|_e = H_{e,\eta}g$  for all edges  $e \subset \partial\Omega$ .

These operators provide stable decompositions of  $L^2(\Omega)$  and  $L^2(\partial\Omega)$ . For example, one has  $L_{\Omega,\eta} + H_{\Omega,\eta} = \text{Id}$  on  $L^2(\Omega)$  and the bounds

$$\|L_{\Omega,\eta}f\|_{L^2(\Omega)} + \|H_{\Omega,\eta}f\|_{L^2(\Omega)} \leq C\|f\|_{L^2(\Omega)} \quad \forall f \in L^2(\Omega),$$

where  $C > 0$  depends solely on  $\Omega$  (via the Stein extension operator  $E_\Omega$ ). The operators  $H_{\Omega,\eta}$  and  $H_{\partial\Omega,\eta}$  have furthermore approximation properties if the function they are applied to has some Sobolev regularity. We illustrate this for  $H_{\partial\Omega,\eta}$ :

**Lemma 3.4.** *Let  $\Omega \subset \mathbb{R}^2$  be a polygon. Then there exists  $C > 0$  independent of  $k$  and  $\eta > 1$  such that for all  $g \in H_{pw}^{1/2}(\partial\Omega)$*

$$k^{1/2}(1 + \eta^{1/2})\|H_{\partial\Omega,\eta}g\|_{L^2(\partial\Omega)} + \|H_{\partial\Omega,\eta}g\|_{H_{pw}^{1/2}(\partial\Omega)} \leq C\|g\|_{H_{pw}^{1/2}(\partial\Omega)}.$$

*Proof.* We only show the estimate for  $\|H_{\partial\Omega,\eta}g\|_{L^2(\partial\Omega)}$ . We consider first the case of an interval  $I \subset \mathbb{R}$ . We define  $H_{I,\eta}g$  by  $H_{I,\eta}g = \mathcal{F}^{-1}\chi_{\mathbb{R}\setminus B_{\eta k}(0)}\mathcal{F}E_I g$ , where  $\chi_{\mathbb{R}\setminus B_{\eta k}(0)}$  is the characteristic function for  $\mathbb{R} \setminus (-\eta k, \eta k)$  and  $E_I$  is the Stein extension operator for the interval  $I$ . Since, by Parseval,  $\mathcal{F}$  is an isometry on  $L^2(\mathbb{R})$  we have

$$\begin{aligned} \|H_{I,\eta}g\|_{L^2(I)}^2 &\leq \|H_{I,\eta}g\|_{L^2(\mathbb{R})}^2 = \int_{\mathbb{R}\setminus B_{\eta k}(0)} |\mathcal{F}E_I g|^2 d\xi \\ &= \int_{\mathbb{R}\setminus B_{\eta k}(0)} \frac{(1 + |\xi|^2)^{1/2}}{(1 + |\xi|^2)^{1/2}} |\mathcal{F}E_I g|^2 d\xi \leq \frac{1}{(1 + (\eta k)^2)^{1/2}} \int_{\mathbb{R}} (1 + |\xi|^2)^{1/2} |\mathcal{F}E_I g|^2 d\xi. \end{aligned}$$

The last integral can be bounded by  $C\|E_I g\|_{H^{1/2}(\mathbb{R})}^2$ . The stability properties of the extension operator  $E_I$  then imply furthermore  $\|E_I g\|_{H^{1/2}(\mathbb{R})} \leq C\|g\|_{H^{1/2}(I)}$ . In total,

$$\|H_{I,\eta}g\|_{L^2(I)} \leq C \frac{1}{(1 + (\eta k)^2)^{1/4}} \|g\|_{H^{1/2}(I)} \leq Ck^{-1/2}(1 + \eta)^{-1/2}\|g\|_{H^{1/2}(I)},$$

where, in the last estimate, the constant  $C$  depends additionally on  $k_0$ . From this estimate, we obtain the desired bound for  $\|H_{\partial\Omega,\eta}g\|_{L^2(\partial\Omega)}$  by identifying each edge of  $\Omega$  with an interval.  $\square$

### 3.3.2 Corner singularities

We recall the following result harking back to the work by Kondratiev and Grisvard:

**Lemma 3.5.** *Let  $\Omega \subset \mathbb{R}^d$  be a polygon with vertices  $A_j$ ,  $j = 1, \dots, J$ , and interior angles  $\omega_j$ ,  $j = 1, \dots, J$ . Define for each vertex  $A_j$  the singularity function  $S_j$  by*

$$S_j(r_j, \varphi_j) = r_j^{\pi/\omega_j} \cos\left(\frac{\pi}{\omega_j} \varphi_j\right), \quad (26)$$

where  $(r_j, \varphi_j)$  are polar coordinates centered at the vertex  $A_j$  such that the edges of  $\Omega$  meeting at  $A_j$  correspond to  $\varphi_j = 0$  and  $\varphi_j = \omega_j$ . Then every solution  $u$  of

$$-\Delta u = f \quad \text{in } \Omega, \quad \partial_n u = g \quad \text{on } \partial\Omega,$$

can be written as  $u = u_0 + \sum_{j=1}^J a_j^A(f, g) S_j$  with the a priori bounds

$$\|u_0\|_{H^2(\Omega)} + \sum_{j=1}^J |a_j^A(f, g)| \leq C \left[ \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)} + \|u\|_{H^1(\Omega)} \right]. \quad (27)$$

The  $a_j^A$  are linear functionals, and  $a_j^A = 0$  for convex corners  $A_j$  (i.e., if  $\omega_j < \pi$ ).

*Proof.* This classical result is comprehensively treated in [38].  $\square$

### 3.3.3 The modified Helmholtz equation

We consider the modified Helmholtz equation in both a bounded domain with Robin boundary conditions and in the full space  $\mathbb{R}^2$ . The corresponding solution operators will be denoted  $S_\Omega^+$  and  $S_{\mathbb{R}^2}^+$ :

1. The operator  $S_\Omega^+ : L^2(\Omega) \times H_{pw}^{1/2}(\partial\Omega) \rightarrow H^1(\Omega)$  is the solution operator for

$$-\Delta u + k^2 u = f \quad \text{in } \Omega, \quad \partial_n u + iku = g \quad \text{on } \partial\Omega. \quad (28)$$

2. The operator  $S_{\mathbb{R}^2}^+ : L^2(\mathbb{R}^2) \rightarrow H^1(\mathbb{R}^2)$  is the solution operator for

$$-\Delta u + k^2 u = f \quad \text{in } \mathbb{R}^2. \quad (29)$$

**Lemma 3.6 (properties of  $S_\Omega^+$ ).** *Let  $\Omega \subset \mathbb{R}^2$  be a polygon and  $f \in L^2(\Omega)$ ,  $g \in H_{pw}^{1/2}(\partial\Omega)$ . Then the solution  $u := S_\Omega^+(f, g)$  satisfies*

$$\|u\|_{1,k,\Omega} \leq k^{-1/2} \|g\|_{L^2(\partial\Omega)} + k^{-1} \|f\|_{L^2(\Omega)}. \quad (30)$$

Furthermore, there exists  $C > 0$  independent of  $k$  and the data  $f, g$ , and there exists a decomposition  $u = u_{H^2} + \sum_{i=1}^J a_i^+(f, g) S_i$  for some linear functionals  $a_i^+$  with

$$\|u_{H^2}\|_{H^2(\Omega)} + \sum_{i=1}^J |a_i^+(f, g)| \leq C \left[ \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)} + k^{1/2} \|g\|_{L^2(\partial\Omega)} \right]. \quad (31)$$



*Proof.* The estimate (30) for  $\|u\|_{1,k,\Omega}$  follows by Lax-Milgram – see [62, Lemma 4.6] for details. Since  $u$  satisfies

$$-\Delta u = f - k^2 u \quad \text{in } \Omega, \quad \partial_n u = g - \mathbf{i}ku \quad \text{on } \partial\Omega,$$

the standard regularity theory for the Laplacian (see Lemma 3.5) permits us to decompose  $u = u_{H^2} + \sum_{i=1}^J a_i^\Delta (f - k^2 u, g - \mathbf{i}ku) S_i$ . The continuity of the linear functionals  $a_i^\Delta$  reads

$$\sum_{i=1}^J |a_i^\Delta (f - k^2 u, g - \mathbf{i}ku)| \leq C \left[ \|f - k^2 u\|_{L^2(\Omega)} + \|g - \mathbf{i}ku\|_{H_{pw}^{1/2}(\partial\Omega)} \right].$$

Since  $(f, g) \mapsto S_\Omega^+(f, g)$  is linear, the map  $(f, g) \mapsto a_i^+(f, g) := a_i^\Delta (f - k^2 u, g - \mathbf{i}ku)$  is linear, and (30), (27) give the desired estimates for  $u_{H^2}$  and  $a_i^+(f, g)$ .  $\square$

**Lemma 3.7 (properties of  $S_{\mathbb{R}^2}^+$ ).** *There exists  $C > 0$  such that for every  $\eta > 1$  and every  $f \in L^2(\mathbb{R}^2)$  whose Fourier transform  $\mathcal{F}f$  satisfies  $\text{supp } \mathcal{F}f \subset \mathbb{R}^2 \setminus B_{\eta k}(0)$ , the solution  $u = S_{\mathbb{R}^2}^+ f$  of (29) satisfies*

$$\|u\|_{1,k,\mathbb{R}^2} \leq k^{-1} \frac{1}{\sqrt{1+\eta^2}} \|f\|_{L^2(\mathbb{R}^2)}, \quad \|u\|_{H^2(\mathbb{R}^2)} \leq C \|f\|_{L^2(\mathbb{R}^2)}.$$

*Proof.* The result follows from Parseval's theorem and the weak formulation for  $u$  as follows (we abbreviate the Fourier transforms by  $\widehat{f} = \mathcal{F}f$  and  $\widehat{u} = \mathcal{F}u$ ):

$$\begin{aligned} \|u\|_{1,k,\mathbb{R}^2}^2 &= (f, u)_{L^2(\mathbb{R}^2)} = (\widehat{f}, \widehat{u})_{L^2(\mathbb{R}^2)} \\ &\leq \sqrt{\int_{\mathbb{R}^2} (|\xi|^2 + k^2)^{-1} |\widehat{f}|^2 d\xi} \sqrt{\int_{\mathbb{R}^2} (|\xi|^2 + k^2) |\widehat{u}|^2 d\xi} \\ &= \sqrt{\int_{\mathbb{R}^2 \setminus B_{\eta k}(0)} (|\xi|^2 + k^2)^{-1} |\widehat{f}|^2 d\xi} \|u\|_{1,k,\mathbb{R}^2} \leq \frac{1}{k\sqrt{1+\eta^2}} \|\widehat{f}\|_{L^2(\mathbb{R}^2)} \|u\|_{1,k,\mathbb{R}^2}, \end{aligned}$$

where, in the penultimate step, we used the support properties of  $\widehat{f}$ . Appealing again to Parseval, we get the desired claim for  $\|u\|_{1,k,\mathbb{R}^2}$ . The estimate for  $\|u\|_{H^2(\mathbb{R}^2)}$  now follows from  $f \in L^2(\mathbb{R}^2)$  and the standard interior regularity for the Laplacian.  $\square$

### 3.4 Proof of Theorem 3.2

We denote by  $S : (f, g) \mapsto S(f, g)$  the solution operator for (9). Concerning some of its properties, we have the following result taken essentially from [62, Lemma 4.13]:

**Lemma 3.8 (analytic regularity of  $S(f, g)$ ).** *Let  $\Omega$  be a polygon. Let  $f$  be analytic on  $\Omega$  and  $g \in L^2(\partial\Omega)$  be piecewise analytic and satisfy for some constants  $\widetilde{C}_f, \widetilde{C}_g$ ,*

$\gamma_f, \gamma_g > 0$

$$\|\nabla^n f\|_{L^2(\Omega)} \leq \tilde{C}_f \gamma_f^n \max\{n, k\}^n \quad \forall n \in \mathbb{N}_0 \quad (32a)$$

$$\|\nabla_T^n g\|_{L^2(e)} \leq \tilde{C}_g \gamma_g^n \max\{n, k\}^n \quad \forall n \in \mathbb{N}_0 \quad \forall e \in \mathcal{E}, \quad (32b)$$

where  $\mathcal{E}$  denotes the set of edges of  $\Omega$  and  $\nabla_T$  tangential differentiation. Then there exist  $\vec{\beta} \in [0, 1]^J$  (depending only on  $\Omega$ ) and constants  $C, \gamma > 0$  (depending only on  $\Omega, \gamma_f, \gamma_g, k_0$ ) such that the following is true with the constant  $C_{sol}(k)$  of (25):

$$\|u\|_{1,k,\Omega} \leq C_{sol}(k)(\tilde{C}_f + \tilde{C}_g) \quad (33)$$

$$\|\phi_{n,\vec{\beta},k} \nabla^{n+2} u\|_{L^2(\Omega)} \leq CC_{sol}(k)k^{-1}(\tilde{C}_f + \tilde{C}_g)\gamma^n \max\{n, k\}^{n+2} \quad \forall n \in \mathbb{N}_0. \quad (34)$$

*Proof.* The estimate (33) is simply a restatement of the definition of  $C_{sol}(k)$ . The estimate (34) will follow from [57, Prop. 5.4.5]. To simplify the presentation, we assume by linearity that  $g$  vanishes on all edges of  $\Omega$  with the exception of one edge  $\Gamma$ . Furthermore, we restrict our attention to the vicinity of one vertex, which we take to be the origin; we assume  $\Gamma \subset (0, \infty) \times \{0\}$ , and that near the origin,  $\Omega$  is above  $(0, \infty) \times \{0\}$ , i.e.,  $\{(r \cos \varphi, r \sin \varphi) : 0 < r < \rho, 0 < \varphi < \omega\} \subset \Omega$  for some  $\rho, \omega > 0$ .

Upon setting  $\varepsilon := 1/k$ , we note that  $u$  solves

$$-\varepsilon^2 \Delta u - u = \varepsilon^2 f \quad \text{on } \Omega, \quad \varepsilon^2 \partial_n u = \varepsilon(\varepsilon g - iu) \quad \text{on } \partial\Omega.$$

On the edge  $\Gamma$ , the function  $g$  is the restriction of  $G_{1,0}(x, y) := g(x)e^{-y/\varepsilon}$  to  $\Gamma$ . The assumptions on  $f$  and  $g$  then imply that [57, Prop. 5.4.5] is applicable with the following choice of constants appearing in [57, Prop. 5.4.5]:

$$C_f = \varepsilon^2 \tilde{C}_f, \quad C_{G_1} = \varepsilon \varepsilon^{1/2} \tilde{C}_g, \quad C_{G_2} = \varepsilon, \quad C_b = 0, \quad C_c = 1, \\ \gamma_f = O(1), \quad \gamma_{G_1} = O(1), \quad \gamma_{G_2} = O(1), \quad \gamma_b = 0, \quad \gamma_c = 0,$$

resulting in the existence of constants  $C, K > 0$  and  $\vec{\beta} \in [0, 1]^J$  with

$$\|\Phi_{n,\vec{\beta},k} \nabla^{n+2} u\|_{L^2(\Omega)} \leq CK^{n+2} \max\{n+2, k\}^{n+2} \left( k^{-2} \tilde{C}_f + k^{-1} \|u\|_{1,k,\Omega} + k^{-3/2} \tilde{C}_g \right)$$

for all  $n \in \mathbb{N}_0$ . We conclude the argument by inserting (33) and estimating generously  $k^{-1} \tilde{C}_f + k^{-1/2} \tilde{C}_g \leq C(\tilde{C}_f + \tilde{C}_g)$ .

We remark that this last generous estimate comes from the precise form of our stability assumption (25). Its form (25) is motivated by the estimates *available* for the star-shaped case, but could clearly be replaced with other assumptions.  $\square$

**Corollary 3.9 (analytic regularity of  $S(L_{\Omega,\eta}f, L_{\partial\Omega,\eta}g)$ ).** *Let  $\Omega$  be a polygon and  $\eta > 1$ . Then there exist  $\vec{\beta} \in [0, 1]^J$  (depending only on  $\Omega$ ) and  $C, \gamma > 0$  (depending only on  $\Omega, k_0$ , and  $\eta > 1$ ) such that for every  $f \in L^2(\Omega)$  and  $g \in L^2(\partial\Omega)$ , the function  $u = S(L_{\Omega,\eta}f, L_{\partial\Omega,\eta}g)$  satisfies with  $C_{f,g} := \|f\|_{L^2(\Omega)} + \|g\|_{L^2(\partial\Omega)}$*

$$\|u\|_{1,k,\Omega} \leq CC_{sol}(k)C_{f,g} \quad (35)$$

$$\|\Phi_{n,\vec{\beta},k} \nabla^{n+2} u\|_{L^2(\Omega)} \leq CC_{sol}(k)k^{-1}\gamma^n \max\{n,k\}^{n+2}C_{f,g} \quad \forall n \in \mathbb{N}_0. \quad (36)$$

*Proof.* The definitions of  $L_{\Omega,\eta}f$  and  $L_{\partial\Omega,\eta}g$  imply with Parseval

$$\begin{aligned} \|\nabla^n L_{\Omega,\eta}f\|_{L^2(\Omega)} &\leq C(\eta k)^n \|f\|_{L^2(\Omega)} \quad \forall n \in \mathbb{N}_0, \\ \|\nabla_T^n L_{\partial\Omega,\eta}g\|_{L^2(\partial\Omega)} &\leq C(\eta k)^n \|g\|_{L^2(\partial\Omega)} \quad \forall n \in \mathbb{N}_0, \end{aligned}$$

where again  $\nabla_T$  is the (edgewise) tangential gradient. The desired estimates now follow from Lemma 3.8.  $\square$

Key to the proof of Theorem 3.2 is the following contraction result:

**Lemma 3.10 (contraction lemma).** *Let  $\Omega \subset \mathbb{R}^2$  be a polygon. Fix  $q \in (0, 1)$ . Then one can find  $\vec{\beta} \in [0, 1]^J$  (depending solely on  $\Omega$ ) and constants  $C, \gamma > 0$  independent of  $k$  such that for every  $f \in L^2(\Omega)$  and every  $g \in H_{pw}^{1/2}(\partial\Omega)$ , the solution  $u$  of (9) can be decomposed as  $u = u_{H^2} + \sum_{i=1}^J a_i(f, g)S_i + u_{\mathcal{A}} + r$ , where  $u_{H^2} \in H^2(\Omega)$ , the  $a_i$  are linear functionals, and  $u_{\mathcal{A}} \in C^\infty(\Omega)$ . These functions satisfy*

$$\begin{aligned} k\|u_{H^2}\|_{1,k,\Omega} + \|u_{H^2}\|_{H^2(\Omega)} + \sum_{i=1}^J |a_i(f, g)| &\leq C \left[ \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)} \right], \\ \|u_{\mathcal{A}}\|_{1,k,\Omega} &\leq CC_{sol}(k) \left[ \|f\|_{L^2(\Omega)} + \|g\|_{L^2(\partial\Omega)} \right], \\ \|\Phi_{n,\vec{\beta},k} \nabla^{n+2} u_{\mathcal{A}}\|_{L^2(\Omega)} &\leq CC_{sol}(k)k^{-1}\gamma^n \max\{n,k\}^{n+2} \left[ \|f\|_{L^2(\Omega)} + \|g\|_{L^2(\partial\Omega)} \right] \end{aligned}$$

for all  $n \in \mathbb{N}_0$ . Finally, the remainder  $r$  satisfies

$$-\Delta r - k^2 r = \tilde{f} \quad \text{on } \Omega, \quad \partial_n r + \mathbf{i}kr = \tilde{g}$$

for some  $\tilde{f} \in L^2(\Omega)$  and  $\tilde{g} \in H_{pw}^{1/2}(\partial\Omega)$  with

$$\|\tilde{f}\|_{L^2(\Omega)} + \|\tilde{g}\|_{H_{pw}^{1/2}(\partial\Omega)} \leq q \left( \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)} \right).$$

*Proof.* We start by decomposing  $(f, g) = (L_{\Omega,\eta}f, L_{\partial\Omega,\eta}g) + (H_{\Omega,\eta}f, H_{\partial\Omega,\eta}g)$  with a parameter  $\eta > 1$  that will be selected below. We set

$$u_{\mathcal{A}} := S(L_{\Omega,\eta}f, L_{\partial\Omega,\eta}g), \quad u_1 := S_{\mathbb{R}^2}^+(H_{\Omega,\eta}f),$$

where we tacitly extended  $H_{\Omega,\eta}f$  (which is only defined on  $\Omega$ ) by zero outside  $\Omega$ . Then  $u_{\mathcal{A}}$  satisfies the desired estimates by Corollary 3.9. For  $u_1$  we have by Lemma 3.7 and the stability  $\|H_{\Omega,\eta}f\|_{L^2(\Omega)} \leq C\|f\|_{L^2(\Omega)}$  (we note that  $C > 0$  is independent of  $k$  and  $\eta$ ) the *a priori* estimates

$$\begin{aligned} \|u_1\|_{1,k,\mathbb{R}^2} &\leq Ck^{-1}(1+\eta^2)^{-1/2}\|H_{\Omega,\eta}f\|_{L^2(\Omega)} \leq Ck^{-1}(1+\eta)^{-1}\|f\|_{L^2(\Omega)}, \\ \|u_1\|_{H^2(\mathbb{R}^2)} &\leq C\|H_{\Omega,\eta}f\|_{L^2(\Omega)} \leq C\|f\|_{L^2(\Omega)}. \end{aligned}$$

The trace and the multiplicative trace inequalities imply for  $g_1 := \partial_n u_1 + \mathbf{i}k u_1$ :

$$k^{1/2}(1+\eta)^{1/2}\|g_1\|_{L^2(\partial\Omega)} + \|g_1\|_{H_{pw}^{1/2}(\partial\Omega)} \leq C\|f\|_{L^2(\Omega)}.$$

For  $g_2 := H_{\partial\Omega,\eta}g - g_1$  we then get from Lemma 3.4 and the triangle inequality

$$k^{1/2}(1+\eta)^{1/2}\|g_2\|_{L^2(\partial\Omega)} + \|g_2\|_{H_{pw}^{1/2}(\partial\Omega)} \leq C \left[ \|g\|_{H_{pw}^{1/2}(\partial\Omega)} + \|f\|_{L^2(\Omega)} \right].$$

Lemma 3.6 yields for  $u_2 := S_{\Omega}^+(0, g_2)$ ,

$$\|u_2\|_{1,k,\Omega} \leq Ck^{-1/2}\|g_2\|_{L^2(\partial\Omega)} \leq Ck^{-1}(1+\eta)^{-1/2} \left[ \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)} \right],$$

and furthermore we can write  $u_2 = u_{H^2} + \sum_{i=1}^J a_i^+(0, g_2) S_i$ , with

$$\|u_{H^2}\|_{H^2(\Omega)} + \sum_{i=1}^J |a_i^+(0, g_2)| \leq C \left[ \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)} \right].$$

We then define  $a_i(f, g) := a_i^+(0, g_2)$  and note that  $(f, g) \mapsto a_i(f, g)$  is linear by linearity of the maps  $a_i^+$  and  $(f, g) \mapsto g_2$ . The above shows that  $u_{H^2}$  and the  $a_i$  satisfy the required estimates. Finally, the function  $\tilde{u} := u - (u_{\mathcal{A}} + u_1 + u_2)$  satisfies

$$-\Delta \tilde{u} - k^2 \tilde{u} = 2k^2(u_1 + u_2) =: \tilde{f}, \quad \partial_n \tilde{u} + \mathbf{i}k \tilde{u} = 0 =: \tilde{g},$$

together with

$$\|\tilde{f}\|_{L^2(\Omega)} \leq 2k^2 \left( \|u_1\|_{L^2(\Omega)} + \|u_2\|_{L^2(\Omega)} \right) \leq C(1+\eta)^{-1/2} \left[ \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)} \right].$$

Hence, selecting  $\eta > 1$  sufficiently large so that for the chosen  $q \in (0, 1)$  we have  $C(1+\eta)^{-1/2} \leq q$  allows us to conclude the proof.  $\square$

*Proof of Theorem 3.2.* The contraction property of Lemma 3.10 can be restated as  $S(f, g) = u_{H^2} + \sum_{i=1}^J a_i(f, g) S_i + u_{\mathcal{A}} + S(\tilde{f}, \tilde{g})$ , where, for a chosen  $q \in (0, 1)$ , we have  $\|\tilde{f}\|_{L^2(\Omega)} + \|\tilde{g}\|_{H_{pw}^{1/2}(\partial\Omega)} \leq q \left[ \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)} \right]$  together with appropriate estimates for  $u_{H^2}$ ,  $a_i(f, g)$ , and  $u_{\mathcal{A}}$ . This consideration can be repeated for  $S(\tilde{f}, \tilde{g})$ . We conclude that a geometric series argument can be employed to write  $u = S(f, g) = u_{H^2} + \sum_{i=1}^J \tilde{a}_i(f, g) S_i + \tilde{u}_{\mathcal{A}}$ , where  $u_{H^2} \in H^2(\Omega)$ ,  $\tilde{u}_{\mathcal{A}} \in C^\infty(\Omega)$ , and the coefficients  $\tilde{a}_i$  are in fact linear functionals of the data  $(f, g)$ . Furthermore, we have with the abbreviation  $C_{f,g} := \|f\|_{L^2(\Omega)} + \|g\|_{H_{pw}^{1/2}(\partial\Omega)}$

$$\begin{aligned}
\|\tilde{u}_{\mathcal{A}}\|_{1,k,\Omega} &\leq CC_{f,g} \\
\|\Phi_{n,\vec{\beta},k} \nabla^{n+2} \tilde{u}_{\mathcal{A}}\|_{L^2(\Omega)} &\leq CC_{sol}(k)k^{-1}C_{f,g}\gamma^n \max\{n,k\}^{n+2} \quad \forall n \in \mathbb{N}_0, \\
k\|u_{H^2}\|_{1,k,\Omega} + \|u_{H^2}\|_{H^2(\Omega)} + \sum_{i=1}^J |\tilde{a}_i(f,g)| &\leq CC_{f,g}.
\end{aligned}$$

Finally, Lemma 3.11 below allows us to absorb the contribution  $\sum_{i=1}^J \tilde{a}_i(f,g)S_i$  in the analytic part by setting  $u_{\mathcal{A}} := \tilde{u}_{\mathcal{A}} + \sum_{i=1}^J \tilde{a}_i(f,g)S_i$ . In view of  $\beta_i < 1$ , we have  $2 - \beta_i \geq 1$  and arrive at

$$\begin{aligned}
\|u_{\mathcal{A}}\|_{H^1(\Omega)} &\leq C(C_{sol}(k) + 1)C_{f,g}, \quad k\|u_{\mathcal{A}}\|_{L^2(\Omega)} \leq CC_{f,g}(C_{sol}(k) + k), \\
\|\Phi_{n,\vec{\beta},k} \nabla^{n+2} u_{\mathcal{A}}\|_{L^2(\Omega)} &\leq CC_{f,g} [C_{sol}(k)k^{-1} + k^{-1}] \max\{n,k\}^{n+2} \quad \forall n \in \mathbb{N}_0,
\end{aligned}$$

which concludes the argument.  $\square$

**Lemma 3.11.** *Let  $\beta_i \in [0, 1)$  satisfy  $\beta_i > 1 - \frac{\pi}{\omega_i}$ . Then, for some  $C, \gamma > 0$  independent of  $k$ , the singularity functions  $S_i$  of (26) satisfy  $\|S_i\|_{H^1(\Omega)} \leq C$  and*

$$\|\Phi_{n,\vec{\beta},k} \nabla^{n+2} S_i\|_{L^2(\Omega)} \leq Ck^{-(2-\beta_i)}\gamma^n \max\{n,k\}^{n+2} \quad \forall n \in \mathbb{N}_0$$

*Proof.* Follows by a direct calculation. See Lemma A.3 for details.  $\square$

## 4 Stability of Galerkin discretizations

### 4.1 Abstract results

We consider the model problem (9) and a sequence  $(V_N)_{N \in \mathbb{N}} \subset H^1(\Omega)$  of finite-dimensional spaces. Furthermore, we assume that  $(V_N)_{N \in \mathbb{N}}$  is such that for every  $v \in H^1(\Omega)$  we have  $\lim_{N \rightarrow \infty} \inf_{v_N \in V_N} \|v - v_N\|_{H^1(\Omega)} = 0$ . The conforming approximations  $u_N$  to the solution  $u$  of (9) are then defined by:

$$\text{Find } u_N \in V_N \text{ s.t. } B(u_N, v) = l(v) \quad \forall v \in V_N. \quad (37)$$

Since the sesquilinear form  $B$  satisfies a Gårding inequality, general functional analytic argument show that *asymptotically*, the discrete problem (37) has a unique solution  $u_N$  and are quasi-optimal (see, e.g., [73, Thm. 4.2.9], [74]). More precisely, there exist  $N_0 > 0$  and  $C > 0$  such that for all  $N \geq N_0$

$$\|u - u_N\|_{1,k,\Omega} \leq C \inf_{v \in V_N} \|u - v\|_{1,k,\Omega}. \quad (38)$$

This general functional analytic argument does not give any indication of how  $C$  and  $N_0$  depend on discretization parameters and the wavenumber  $k$ . Inspection of the arguments reveals that it is the approximation properties of the spaces  $V_N$  for the

approximation of the solution of certain adjoint problems that leads to the quasi-optimality result (38). For the reader's convenience, we repeat the argument, which has been used previously in, e.g., [6, 13, 56, 61, 62, 72, 74] and is often attributed to Schatz, [74]:

**Lemma 4.1 ([62, Thm. 3.2]).** *Let  $\Omega \subset \mathbb{R}^d$  be a bounded Lipschitz domain and  $B$  be defined in (11). Denote by  $S^* : L^2(\Omega) \rightarrow H^1(\Omega)$  the solution operator for the problem*

$$\text{Find } u^* \in H^1(\Omega) \text{ s.t. } B(v, u^*) = (v, f)_{L^2(\Omega)} \quad \forall v \in H^1(\Omega). \quad (39)$$

Define the adjoint approximation property  $\eta(V_N)$  by

$$\eta(V_N) := \sup_{f \in L^2(\Omega)} \inf_{v \in V_N} \frac{\|S^*(f) - v\|_{1,k,\Omega}}{\|f\|_{L^2(\Omega)}}.$$

If, for the continuity constant  $C_B$  of (12), the space  $V_N$  satisfies

$$2C_B k \eta(V_N) \leq 1, \quad (40)$$

then the solution  $u_N$  of (37) exists and satisfies

$$\|u - u_N\|_{1,k,\Omega} \leq 2C_B \inf_{v \in V_N} \|u - v\|_{1,k,\Omega}. \quad (41)$$

*Proof.* We will not show existence of  $u_N$  but restrict our attention on the quasi-optimality result (41); we refer to [54, Thm. 3.9] for the demonstration that (41) in fact implies existence and uniqueness of  $u_N$ . Letting  $e = u - u_N$  be the error, we start with an estimate for  $\|e\|_{L^2(\Omega)}$ : Using the definition of the operator  $S^*$  and the Galerkin orthogonality satisfied by  $e$ , we have for arbitrary  $v \in V_N$

$$\|e\|_{L^2(\Omega)}^2 = (e, e)_{L^2(\Omega)} = B(e, S^*e) = B(e, S^*e - v) \leq C_B \|e\|_{1,k,\Omega} \|S^*e - v\|_{1,k,\Omega}.$$

Infimizing over all  $v \in V_N$  yields with the adjoint approximation property  $\eta(V_N)$

$$\|e\|_{L^2(\Omega)} \leq C_B \eta(V_N) \|e\|_{1,k,\Omega}.$$

The Gårding inequality and the Galerkin orthogonality yield for arbitrary  $v \in V_N$ :

$$\begin{aligned} \|e\|_{1,k,\Omega}^2 &= \operatorname{Re} B(e, e) + 2k^2 \|e\|_{L^2(\Omega)}^2 = \operatorname{Re} B(e, u - v) + 2k^2 \|e\|_{L^2(\Omega)}^2 \\ &\leq C_B \|e\|_{1,k,\Omega} \|u - v\|_{1,k,\Omega} + (C_B k \eta(V_N))^2 \|e\|_{1,k,\Omega}^2. \end{aligned}$$

The assumption  $C_B k \eta(V_N) \leq 1/2$  allows us to rearrange this bound to get  $\|e\|_{1,k,\Omega} \leq 2C_B \|u - v\|_{1,k,\Omega}$ . Since  $v \in V_N$  is arbitrary, we arrive at (41).  $\square$

Lemma 4.1 informs us that the convergence analysis for the Galerkin discretization of (9) can be reduced to the study of the adjoint approximation property  $\eta(V_N)$ , which is purely a question of approximation theory. In the context of piecewise

polynomial approximation spaces  $V_N$  this requires a good regularity theory for the operator  $S^*$ . The strong form of the equation satisfied by  $u^* := S^*f$  is

$$-\Delta u^* - k^2 u^* = f \quad \text{in } \Omega, \quad \partial_n u^* - iku^* = 0 \quad \text{on } \partial\Omega, \quad (42)$$

which is again a Helmholtz problem of the type considered in Section 3. More formally, with the solution operator  $S$  of Section 3, we have  $S^*f = \overline{S(\overline{f}, 0)}$ , where an overbar denotes complex conjugation. Thus, the regularity theory of Section 3 is applicable.

## 4.2 Stability of $hp$ -FEM

The estimates of Theorem 3.2 suggest that the effect of the corner singularities is essentially restricted to an  $O(1/k)$ -neighborhood of the vertices. This motivates us to consider meshes that are refined in a small neighborhood of the vertices. To fix ideas, we restrict our attention to meshes  $\mathcal{T}_{h,L}^{geo}$  that are obtained in the following way: First, a quasi-uniform triangulation  $\mathcal{T}_h$  with mesh size  $h$  is selected. Then, the elements abutting the vertices  $A_j$ ,  $j = 1, \dots, J$ , are refined further with a mesh that is geometrically graded towards these vertices. These geometric meshes have  $L$  layers and use a grading factor  $\sigma \in (0, 1)$  (see [77, Sec. 4.4.1] for a precise formal definition). Furthermore, for any regular, shape-regular mesh  $\mathcal{T}$ , we define

$$S^p(\mathcal{T}) := \{u \in H^1(\Omega) : u|_K \in \mathcal{P}_p \quad \forall K \in \mathcal{T}\}, \quad (43)$$

where  $\mathcal{P}_p$  denotes the space of polynomials of degree  $p$ . We now show that on the geometric meshes  $\mathcal{T}_{h,L}^{geo}$ , stability of the FEM is ensured if the mesh size  $h$  and the polynomial degree  $p$  satisfy the scale resolution condition (6) and, additionally,  $L = O(p)$  layers of geometric refinement are used near the vertices:

**Theorem 4.2 (quasi-optimality of  $hp$ -FEM).** *Let  $\mathcal{T}_{h,L}^{geo}$  denote the geometric meshes on the polygon  $\Omega \subset \mathbb{R}^2$  as described above. Fix  $c_3 > 0$ . Then there are constants  $c_1, c_2 > 0$  depending solely on  $\Omega$  and the shape-regularity of the mesh  $\mathcal{T}_{h,L}^{geo}$  such that the following is true: If  $h, p$ , and  $L$  satisfy the conditions*

$$\frac{kh}{p} \leq c_1 \quad \text{and} \quad p \geq c_2 \log k \quad \text{and} \quad L \geq c_3 p \quad (44)$$

*then the  $hp$ -FEM based on the space  $S^p(\mathcal{T}_{h,L}^{geo})$  has a unique solution  $u_N \in S^p(\mathcal{T}_{h,L}^{geo})$  and*

$$\|u - u_N\|_{1,k,\Omega} \leq 2C_B \inf_{v \in S^p(\mathcal{T}_{h,L}^{geo})} \|u - v\|_{1,k,\Omega}. \quad (45)$$

*Proof.* By Lemma 4.1, we have to estimate  $k\eta(V_N)$  with  $V_N = S^p(\mathcal{T}_{h,L}^{geo})$ . Recalling the definition of  $\eta(V_N)$  we let  $f \in L^2(\Omega)$  and observe that we can decompose  $S^*f = u_{H^2} + u_{\mathcal{A}}$ , where  $u_{H^2}$  and  $u_{\mathcal{A}}$  satisfy the bounds

$$\begin{aligned} \|u_{H^2}\|_{H^2(\Omega)} &\leq C\|f\|_{L^2(\Omega)}, \\ \|\Phi_{n,\vec{\beta},k}^{\rightarrow} \nabla^{n+2} u_{\mathcal{A}}\|_{L^2(\Omega)} &\leq C(C_{sol}(k) + 1)k^{-1}\gamma^n \max\{k, n\}^{n+2}\|f\|_{L^2(\Omega)} \quad \forall n \in \mathbb{N}_0. \end{aligned}$$

Piecewise polynomial approximation on  $\mathcal{T}_{h,L}^{geo}$  as discussed in [62, Prop. 5.6] gives under the assumptions  $kh/p \leq C$  and  $L \geq c_3 p$ : (inspection of the proof of [62, Prop. 5.6] shows that only bounds on the derivatives of order  $\geq 2$  are needed):

$$\begin{aligned} \inf_{v \in V_N} \|u_{H^2} - v\|_{1,k,\Omega} &\leq C \frac{h}{p} \|f\|_{L^2(\Omega)}, \\ \inf_{v \in V_N} \|u_{\mathcal{A}} - v\|_{1,k,\Omega} &\leq C \left[ (kh)^{1-\beta_{max}} e^{ckh-bp} + \left( \frac{kh}{\sigma_0 p} \right)^p \right] (C_{sol}(k) + 1) \|f\|_{L^2(\Omega)}, \end{aligned}$$

where  $\beta_{max} = \max_{j=1,\dots,J} \beta_j < 1$ , and  $C, c, b > 0$  are constants independent of  $h, p$ , and  $k$ . From this, we can easily infer

$$k\eta(V_N) \leq C \left\{ \frac{kh}{p} + k(C_{sol}(k) + 1) \left[ (kh)^{1-\beta_{max}} e^{ckh-bp} + \left( \frac{kh}{\sigma_0 p} \right)^p \right] \right\}.$$

Noting that Theorem 2.4 gives  $C_{sol}(k) = O(k^{5/2})$ , and selecting  $c_1$  sufficiently small as well as  $c_2$  sufficient large allows us to make  $k\eta(V_N)$  so small that the condition (40) in Lemma 4.1 is satisfied.  $\square$

**Corollary 4.3 (exponential convergence on geometric meshes).** *Let  $f$  be analytic on  $\overline{\Omega}$  and  $g$  be piecewise analytic, i.e.,  $f, g$  satisfy (32). Given  $c_3 > 0$ , there exist  $c_1, c_2 > 0$  such that under the scale resolution conditions (44) of Theorem 4.2, the finite element approximation  $u_N \in S^p(\mathcal{T}_{h,L}^{geo})$  exists, and there are constants  $C, b > 0$  independent of  $k$  such that the error  $u - u_N$  satisfies*

$$\|u - u_N\|_{1,k,\Omega} \leq C e^{-bp}.$$

*Proof.* In view of Theorem 4.2, estimating  $\|u - u_N\|_{1,k,\Omega}$  is purely a question of approximability for  $c_1$  sufficiently small and  $c_2$  sufficiently large. Lemma 3.8 gives that the solution  $u = S(f, g)$  satisfies the bounds given there and, as in the proof of Theorem 4.2, we conclude from [62, Prop. 5.6] (more precisely, this follows from its proof)

$$\inf_{v \in V_N} \|u_{\mathcal{A}} - v\|_{1,k,\Omega} \leq C \left[ (kh)^{1-\beta_{max}} e^{ckh-bp} + \left( \frac{kh}{\sigma_0 p} \right)^p \right] (C_{sol}(k) + 1) (\tilde{C}_f + \tilde{C}_g).$$

Theorem 2.4 asserts  $C_{sol}(k) = O(k^{5/2})$ , which implies the result by suitably adjusting  $c_1$  and  $c_2$  if necessary.  $\square$

**Remark 4.4.** 1. The problem size  $N = \dim S^p(\mathcal{T}_{h,L}^{geo})$  is  $N = O((L + h^{-2})p^2)$ . The particular choice of  $L = c_3 p$  layers of geometric refinement, approximation order  $p = c_2 \log k$ , and mesh size  $h = c_1 p/k$  in Theorem 4.2 ensures quasi-



- optimality of the  $hp$ -FEM with problem size  $N = O(k^2)$ , i.e., quasi-optimality can be achieved with a fixed number of degrees of freedom per wavelength.
2. The sparsity pattern of the system matrix is that of the classical  $hp$ -FEM, i.e., each row/column has  $O(p^2)$  non-zero entries. Noting that the scale resolution conditions (6), (44) require  $p = O(\log k)$ , we see that the number of non-zero entries per row/column is not independent of  $k$ . It is worth relating this observation to [12]. It is shown there for a model problem in 2D that *no* 9 point stencil can be found that leads to a pollution-free method.
  3. Any space  $V_N$  that contains  $S^p(\mathcal{T}_{h,L}^{geo})$ , where  $h$ ,  $p$ , and  $L$  satisfy the scale resolution condition (44) also features quasi-optimality.
  4. The factor 2 on the right-hand side of (45) is arbitrary and can be replaced by any number greater than 1.
  5. The stability analysis of Theorem 4.2 requires quite a significant mesh refinement near the vertices, namely,  $L \sim p$ . It is not clear whether this is an artifact of the proof. For a more careful numerical analysis of this issue, more detailed information about the stability properties of the solution operator  $S$  is needed, e.g., estimates for  $\|S(f, g)\|_{1,k,B_{1/k}(A_j)}$ .

### 4.3 Numerical examples: $hp$ -FEM

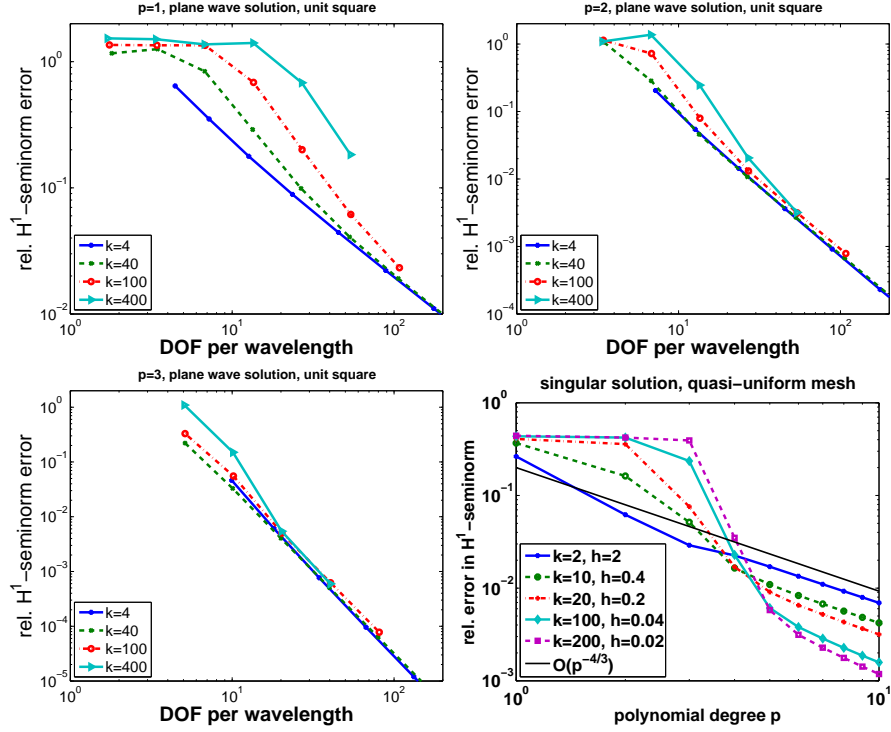
All calculations reported in this section are performed with the  $hp$ -FEM code NETGEN/NGSOLVE by J. Schöberl, [75, 76].

**Example 4.5.** In this 2D analog of Example 1.1, we consider the model problem (9) with exact solution being a plane wave  $e^{i(k_1 x + k_2 y)}$ , where  $k_1 = -k_2 = \frac{1}{\sqrt{2}}k$  and  $k \in \{4, 40, 100, 400\}$ . For fixed  $p \in \{1, 2, 3\}$ , we show in Fig. 2 the performance of the  $h$ -version FEM for  $p \in \{1, 2, 3\}$  on quasi-uniform meshes by displaying the relative error in the  $H^1$ -seminorm versus the number of degrees of freedom per wavelength. We observe that higher order methods are less prone to pollution. We note that the meshes are quasi-uniform, i.e., no geometric mesh refinement near the vertices is performed in contrast to the requirements of Theorem 4.2. ■

**Example 4.6.** On the  $L$ -shaped domain  $\Omega = (-1, 1)^2 \setminus (0, 1) \times (-1, 0)$  with  $\Gamma$  being the union of the two edges meeting at  $(0, 0)$ , we consider

$$-\Delta u - k^2 u = 0 \quad \text{in } \Omega, \quad \partial_n u = 0 \quad \text{on } \Gamma, \quad \partial_n u - iku = g \quad \text{on } \partial\Omega \setminus \Gamma, \quad (46)$$

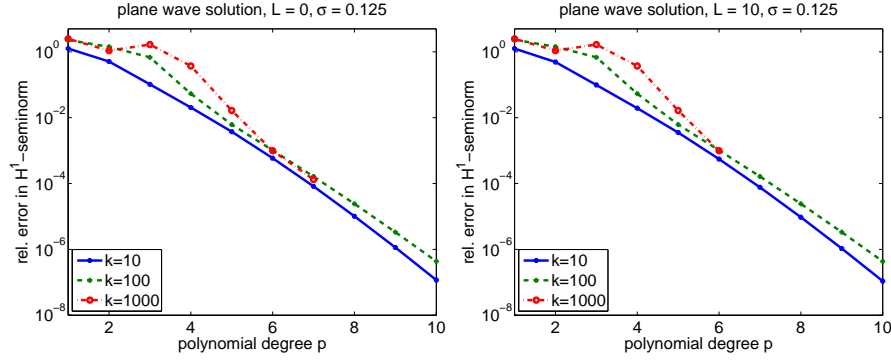
where the Robin data  $g$  are such that the exact solution is  $u(x, y) = e^{i(k_1 x + k_2 y)}$  with  $k_1 = -k_2 = \frac{1}{\sqrt{2}}k$ , and  $k \in \{10, 100, 1000\}$ . We consider two kinds of meshes, namely, quasi-uniform meshes  $\mathcal{T}_h$  with mesh size  $h$  such that  $kh \approx 4$  and meshes  $\mathcal{T}^{geo}$  that are geometrically refined near the origin. The meshes  $\mathcal{T}^{geo}$  are derived from the quasi-uniform mesh  $\mathcal{T}_h$  by introducing a geometric grading on the elements abutting the origin; the grading factor is  $\sigma = 0.125$  and the number of refinement levels is  $L = 10$ . Fig. 3 shows the relative errors in the  $H^1$ -seminorm for the



**Fig. 2** Top:  $h$ -FEM with  $p = 1$  (left) and  $p = 2$  (right) as described in Example 4.5. Bottom left:  $h$ -FEM with  $p = 3$  as described in Example 4.5. Bottom right:  $p$ -FEM for singular solution on quasi-uniform mesh as described in Example 4.7.

$p$ -version of the FEM where for fixed mesh the approximation order  $p$  ranges from 1 to 10. It is particularly noteworthy that the refinement near the origin has hardly any effect on the convergence behavior of the FEM; this is quite in contrast to the stability result Theorem 4.2, which requires geometric refinement near all vertices of  $\Omega$ . ■

**Example 4.7.** The geometry and the boundary conditions are as described in Example 4.6. The data  $g$  are selected such that the exact solution is  $u = J_{2/3}(kr) \cos \frac{2}{3}\varphi$ , where  $(r, \varphi)$  denote polar coordinates and  $J_\alpha$  is a first kind Bessel function.  $k \in \{1, 10, 20, 100, 200\}$ . Our calculations are  $p$ -FEMs with  $p \in \{1, \dots, 10\}$  on the quasiuniform mesh  $\mathcal{T}_h$  described in Example 4.7. The results are displayed in the bottom right part of Fig. 2. The numerics illustrate that the singularity at the origin is rather weak: we observe that the asymptotic algebraic convergence behavior is  $|u - u_N|_{H^1(\Omega)} \approx C_k p^{-4/3} |u|_{H^1(\Omega)}$ , where the constant  $C_k$  depends favorably on  $k$ . ■



**Fig. 3**  $p$ -FEM for plane wave solution as described in Example 4.6. Left: quasiuniform mesh  $\mathcal{T}_h$  with  $kh \approx 4$ . Right: Mesh  $\mathcal{T}^{geo}$  obtained from  $\mathcal{T}_h$  by strong geometric refinement near origin.

#### 4.4 Stability of Partition of Unity Method/Generalized FEM

The abstract stability result of Lemma 4.1 only assumes certain approximation properties of the spaces  $V_N$ . Particularly in an “ $h$ -version” setting, even non-polynomial, operator-adapted spaces may have sufficient approximation properties to ensure the important condition (40) for stability. We illustrate this effect for the PUM/gFEM, [56, 60] with local approximation spaces consisting of systems of plane waves or generalized harmonic polynomials (see Section 5 below) and the classical FEM shape functions as the partition of unity. The key observation is that for  $h$  sufficiently small, the resulting space has approximation properties similar to the classical (low order) FEM space:

**Lemma 4.8.** *Let  $\mathcal{T}$  be a shape-regular triangulation of the polygon  $\Omega \subset \mathbb{R}^2$ . Let  $h$  be its mesh size; let  $(x_i)_{i=1}^M$  be the nodes of the triangulation and  $(\phi_i)_{i=1}^M$  be the piecewise linear hat functions associated with the nodes  $(x_i)_{i=1}^M$ . Assume  $kh \leq C_1$  for some  $C_1 > 0$ . Let  $V^{master}$  be either the space  $V_{GHP}^p$  with  $p \geq 0$  (see (49) below) or the space  $W_{PW}^p$  with  $p \geq 2$  (see (50) below). Define, for each  $i = 1, \dots, M$ , the local approximation  $V_i$  by  $V_i := \text{span}\{b(x - x_i) : b \in V^{master}\}$ . Then the space  $V_N := \sum_{i=1}^M \phi_i V_i$  has the following approximation property: There exists  $C > 0$  depending only on the shape regularity of  $\mathcal{T}$ , the constant  $C_1$ , and  $V^{master}$  such that*

$$\inf_{v \in V_N} \|u - v\|_{L^2(\Omega)} + h \|u - v\|_{H^1(\Omega)} \leq C \left[ h^2 \|u\|_{H^2(\Omega)} + (kh)^2 \|u\|_{L^2(\Omega)} \right] \quad \forall u \in H^2(\Omega).$$

*Proof.* We first show that each local space  $V_i$  has an element  $\psi_i \in V_i$  with

$$h \|\nabla \psi_i\|_{L^\infty(\tilde{\omega}_i)} + \|1 - \psi_i\|_{L^\infty(\tilde{\omega}_i)} \leq C_{app} (kh)^2 \quad (47)$$

for some  $C_{app} > 0$  independent of  $i$  and  $h$ ; here,  $\tilde{\omega}_i = \text{supp } \phi_i$  has diameter  $O(h)$ . It suffices to show (47) for the set  $V^{master}$ . For the space of generalized harmonic polynomials, this follows from  $J_0(kr) = 1 + O((kr)^2)$ , and for the space

of plane waves, Taylor expansion shows that for  $p = 2m$  ( $m \in \mathbb{N}$ ) the function  $\frac{1}{2} [e^{ikx} + e^{-ikx}] = 1 + O((kh)^2)$  has the desired property whereas for odd  $p = 2m + 1$  ( $m \in \mathbb{N}$ ) the observation

$$e^{ik\omega_0 \cdot x} - \frac{1}{2c} [e^{ik\omega_m \cdot x} + e^{ik\omega_{m+1} \cdot x}] = \left(1 - \frac{1}{c}\right) + O((kh)^2), \quad c = \cos \frac{2\pi m}{2m+1},$$

can be utilized to construct  $\psi_i$ . We recall (see, e.g., [18, Thm. 4.8.7]) that for each  $u \in H^2(\Omega)$  there is a function  $w = \sum_i w(x_i) \phi_i \in S^1(\mathcal{T}_h)$  with

$$\|w\|_{L^2(\Omega)} \leq C\|u\|_{L^2(\Omega)}, \quad \|w\|_{H^1(\Omega)} \leq C\|u\|_{H^1(\Omega)}, \quad (48a)$$

$$\|u - w\|_{L^2(\Omega)} \leq Ch^2\|u\|_{H^2(\Omega)}, \quad \|u - w\|_{H^1(\Omega)} \leq Ch\|u\|_{H^2(\Omega)}. \quad (48b)$$

Upon setting  $v := \sum_i w_i(x_i) \psi_i \phi_i \in V_N$ , we get in view of  $\sum_i \phi_i \equiv 1$  for the error  $u - v = u - \sum_i w_i(x_i) \psi_i \phi_i = (u - w) + \sum_i \phi_i w(x_i)(1 - \psi_i)$ . The estimates (48) imply  $\|u - w\|_{L^2(\Omega)} + h\|u - w\|_{H^1(\Omega)} \leq Ch^2\|u\|_{H^2(\Omega)}$ . For the sum, we have

$$\begin{aligned} \left\| \sum_i w_i(x_i)(1 - \psi_i) \phi_i \right\|_{L^2(\Omega)} &\leq C(kh)^2 \|w\|_{L^2(\Omega)} \leq C(kh)^2 \|u\|_{L^2(\Omega)}, \\ h \|\nabla \sum_i w_i(x_i)(1 - \psi_i) \phi_i\|_{L^2(\Omega)} &\leq C(kh)^2 \|w\|_{L^2(\Omega)} \leq C(kh)^2 \|u\|_{L^2(\Omega)}, \end{aligned}$$

which concludes the proof.  $\square$

**Remark 4.9.** The approximation result of Lemma 4.8 can be generalized in various directions. For example, interpolation arguments allow one to construct, for  $v \in H^{1+\theta}(\Omega)$  with  $\theta \in (0, 1)$  an approximation  $v_{app} \in V_N$  such that  $\|v - v_{app}\|_{L^2(\Omega)} + h\|v - v_{app}\|_{H^1(\Omega)} \leq C_1(h^{1+\theta} + (kh)^2 h^\theta) \|v\|_{H^{1+\theta}(\Omega)} + C_2(kh)^2 \|u\|_{L^2(\Omega)}$ . (We refer the reader to the Appendix for the proof of this results.) Furthermore, the approximation result of Lemma 4.8 can be localized, which is of interest if  $\mathcal{T}$  is not quasi-uniform.  $\blacksquare$

Lemma 4.8 shows that the space  $V_N$ , which is derived from solutions of the homogeneous Helmholtz equation, nevertheless has some approximation power for arbitrary functions with some Sobolev regularity. Hence, the condition (40) can be met for sufficiently small mesh sizes:

**Corollary 4.10 ([56, Prop. 8.2.7]).** *Assume the hypotheses of Lemma 4.8; in particular, let the space  $V_N$  be constructed from systems of plane waves or generalized harmonic polynomials. Assume additionally that  $\Omega$  is a convex polygon. Then there exists  $C > 0$  independent of  $k$  such that for  $k^2 h \leq C$  the Galerkin method for (9) with  $f = 0$  is quasi-optimal, i.e., the solution  $u_N \in V_N$  of (37) exists and satisfies*

$$\|u - u_N\|_{1,k,\Omega} \leq 2C_B \inf_{v \in V_N} \|u - v\|_{1,k,\Omega}.$$

*Proof.* In view of Lemma 4.1, we have to estimate  $\eta(V_N)$ . To that end, we consider (9) with  $f \in L^2(\Omega)$  and  $g = 0$ . In view of the convexity of  $\Omega$ , we have  $C_{sol}(k) = O(1)$

and elliptic regularity then yields for the solution  $u$  of (9)

$$\|u\|_{1,k,\Omega} + k^{-1}\|u\|_{H^2(\Omega)} \leq C\|f\|_{L^2(\Omega)}.$$

This allows us to conclude with Lemma 4.8 that

$$\begin{aligned} \inf_{v \in V_N} \|u - v\|_{1,k,\Omega} &\leq C \left[ (kh^2 + h)\|u\|_{H^2(\Omega)} + (k(kh)^2 + k^2h)\|u\|_{L^2(\Omega)} \right] \\ &\leq C((kh)^2 + kh)\|f\|_{L^2(\Omega)} \leq Ckh(1 + kh)\|f\|_{L^2(\Omega)}. \end{aligned}$$

Hence,  $k\eta(V_N)$  can be made sufficiently small if  $k^2h$  is sufficiently small. We point out that convexity of  $\Omega$  is assumed for convenience—under more stringent conditions on the mesh size  $h$ , quasioptimality holds for general polygons.  $\square$

## 5 Approximation with plane, cylindrical, and spherical waves

Systems of functions that are solutions of a (homogeneous) differential equation are often called “Trefftz systems”. Prominent examples in the context of the Helmholtz equation are, in the two-dimensional setting, “generalized harmonic polynomials” and systems of plane waves given by

$$V_{GHP}^p := \text{span}\{J_n(kr)e^{in\varphi} : -p \leq n \leq p\}, \quad (49)$$

$$W_{PW}^p := \text{span}\{e^{ik\omega_n \cdot (x,y)} : n = 0, \dots, p-1\}, \quad \omega_n = (\cos \frac{2\pi n}{p}, \sin \frac{2\pi n}{p}); \quad (50)$$

here,  $J_n$  is a first kind Bessel function, the functions in  $V_{GHP}^p$  are described in polar coordinates and the functions of  $W_{PW}^p$  in Cartesian coordinates. We point out that analogous systems can be developed in 3D. These functions are solutions of the homogeneous Helmholtz equation. For the approximation of a function  $u$  that satisfies the homogeneous Helmholtz equation on a domain  $\Omega \subset \mathbb{R}^2$ , one may study the “ $p$ -version”, i.e., study how well  $u$  can be approximated from the spaces  $V_{GHP}^p$  or  $W_{PW}^p$  as  $p \rightarrow \infty$ ; alternatively, one may study the “ $h$ -version”, in which, for fixed  $p$ , the approximation properties of the spaces  $V_{GHP}^p$  or  $W_{PW}^p$  are expressed in terms of the diameter  $h = \text{diam } \Omega$  of a domain under consideration. In the way of illustration, we present two types of results:

**Lemma 5.1 ([56]).** *Let  $\Omega \subset \mathbb{R}^2$  be a simply connected domain and  $\Omega' \subset\subset \Omega$  be a compact subset. Let  $u$  solve  $-\Delta u - k^2 u = 0$  on  $\Omega$ . Then there exist constants  $C$ ,  $b > 0$  (possibly depending on  $k$ ) such that for all  $p \geq 2$ :*

$$\inf_{v \in V_{GHP}^p} \|u - v\|_{H^1(\Omega')} \leq Ce^{-bp}, \quad \inf_{v \in W_{PW}^p} \|u - v\|_{H^1(\Omega')} \leq Ce^{-bp/\ln p}.$$

*Proof.* See, e.g., [56] or [58, Thm. 5.3].  $\square$

**Remark 5.2.** Analogs of Lemma 5.1 hold if  $u$  has only some finite Sobolev regularity. Then, the convergence rates are algebraic, [56], [58, Thm. 5.4], [42]. ■

The approximation properties of the spaces  $V_{GHP}^p$  and  $W_{PW}^p$  can be also be studied in an  $h$ -version setting:

**Proposition 5.3** ([42, Thm. 3.2.2]). *Let  $\Omega \subset \mathbb{R}^2$  be a domain with diameter  $h$  and inscribed circle of radius  $\rho h$ . Let  $p = 2\mu + 1$ . Assume  $kh \leq C_1$ . Then there exist  $C_p > 0$  (depending only on  $C_1$ ,  $\rho > 0$ ,  $m$ , and  $p$ ) and  $v \in W_{PW}^{2\mu+1}$  such that*

$$\|u - v\|_{j,k,\Omega,\Sigma} \leq C_p h^{\mu-j+1} \|u\|_{\mu+1,k,\Omega,\Sigma}, \quad 0 \leq j \leq \mu + 1,$$

where  $\|v\|_{j,k,\Omega,\Sigma}^2 = \sum_{m=0}^j k^{2(j-m)} |v|_{H^m(\Omega)}^2$ .

Several comments concerning Proposition 5.3 are in order:

1. The constant  $C_p$  in Proposition 5.3 depends favorably on  $p$ , and its dependence on  $p$  can be found in [42, Thm. 3.2.3].
2. Proposition 5.3 is formulated for the space  $W_{PW}^p$  of plane waves—analogous results are valid for generalized harmonic polynomials, see [42, Thm. 2.2.1] for both the  $h$  and  $hp$ -version.
3. Proposition 5.3 is formulated for the two-dimensional case. Similar results are available in 3D, [42].
4. The approximation properties of plane waves in terms of the element size have previously been studied in slightly different norms in [20].

**Remark 5.4.** Plane waves and generalized harmonic polynomials represent by no means the only operator adapted system used in practice. Especially for polygonal domains, the functions  $J_{n\alpha}(kr) \sin(\alpha n\varphi)$ ,  $n \in \mathbb{N}$ , or  $J_{n\alpha}(kr) \sin(\alpha n\varphi)$ ,  $n \in \mathbb{N}_0$ , for suitable  $\alpha$  can combine good approximation properties with the option to realize homogeneous boundary conditions, [14]. Further possibilities include linear combinations of fundamental solutions or, more generally, discretizations of layer potentials. We refer to [14] for a concrete example. ■

## 6 Stability of Least Squares and DG methods

Discrete stability in Section 4 is obtained from stability of the continuous problem by a perturbation argument. This approach does not seem to work very well if one aims at using approximation spaces that have special features linked to the differential equation under consideration. The reason can be seen from the proof of Lemma 4.1: The adjoint approximation property  $\eta(V_N)$  (which needs to be small) measures how well certain solutions to the inhomogeneous equation can be approximated from the test space. If the ansatz space is based on solutions of the homogeneous equation, then its capabilities to approximate solutions of the inhomogeneous equation are clearly limited. In an  $h$ -version, the situation is not as severe as we have

just seen in Section 4.4 for the PUM/gFEM. In a pure  $p$ -version setting, however, the techniques of Section 4.4 do not seem applicable.

An option is to leave the setting of Galerkin methods and to work with formulations with built-in stability properties. Such approaches can often be understood as minimizing some residual norm, which then provides automatically stability and quasi-optimality (in this residual norm). We will illustrate this procedure here by two examples, namely, Least Squares methods and DG-methods. Our presentation will highlight an issue stemming from this approach, namely, the fact that error estimates in this residual norm do not easily lead to error estimates in more classical norms such as the  $L^2(\Omega)$ -norm.

### 6.1 Some notation for spaces of piecewise smooth functions

Let  $\mathcal{T}$  be a regular, shape-regular triangulation of the polygon  $\Omega \subset \mathbb{R}^2$ . We decompose the set of edges  $\mathcal{E}$  as  $\mathcal{E} = \mathcal{E}_I \cup \mathcal{E}_B$ , where  $\mathcal{E}_I$  is the set of edges in  $\Omega$  and  $\mathcal{E}_B$  consists of the edges on  $\partial\Omega$ . For functions  $u : \Omega \rightarrow \mathbb{R}$  and  $\sigma : \Omega \rightarrow \mathbb{R}^2$  that are smooth on the elements  $K \in \mathcal{T}$ , we define the jumps and averages as it is customary in DG-settings:

- For  $e \in \mathcal{E}_I$ , let  $K_e^+$  and  $K_e^-$  be the two elements sharing  $e$  and denote by  $\mathbf{n}^+$  and  $\mathbf{n}^-$  the normal vectors on  $e$  pointing out of  $K_e^+$  and  $K_e^-$ . Correspondingly, we let  $u^+$ ,  $u^-$  and  $\sigma^+$  and  $\sigma^-$  be traces on  $e$  of  $u$  and  $\sigma$  from  $K_e^+$  and  $K_e^-$ . We set:

$$\begin{aligned} \{u\}_e &:= \frac{1}{2} (u^+ + u^-), & \{\sigma\}_e &:= \frac{1}{2} (\sigma^+ + \sigma^-), \\ [[u]]_e &:= u^+ \mathbf{n}^+ + u^- \mathbf{n}^-, & [[\sigma]]_e &:= \sigma^+ \cdot \mathbf{n}^+ + \sigma^- \cdot \mathbf{n}^-. \end{aligned}$$

- For boundary edges  $e \in \mathcal{E}_B$  we define

$$\{\sigma\}_e := \sigma|_e \quad [[u]]_e := u|_e \mathbf{n}$$

With this notation, one can conveniently rearrange certain sums over edges:

**Lemma 6.1 (“DG magic formula”).** *Let  $v : \Omega \rightarrow \mathbb{R}$  and  $\sigma : \Omega \rightarrow \mathbb{R}^2$  be piecewise smooth on the triangulation  $\mathcal{T}$ . Then:*

$$\sum_{K \in \mathcal{T}} \int_{\partial K} v \sigma \cdot \mathbf{n} = \int_{\mathcal{E}_I} [[v]] \cdot \{\sigma\} + \int_{\mathcal{E}_I} \{v\} \cdot [[\sigma]] + \int_{\mathcal{E}_B} [[v]] \cdot \{\sigma\},$$

where  $\int_{\mathcal{E}_I}$  and  $\int_{\mathcal{E}_B}$  are shorthand notations for the sums of integrals over all edges in  $\mathcal{E}_I$  and  $\mathcal{E}_B$ .

Finally, for piecewise smooth functions,  $\nabla_h$  denotes the piecewise defined gradient.

## 6.2 Stability of least squares methods

Although Least Squares methods could be based on any space of approximation spaces, we will concentrate here on the approximation by piecewise solutions of the homogeneous Helmholtz equation. With varying focus, this is the setting of [14, 53, 65, 70, 80] and references therein. We illustrate the procedure for the model problem (9) with  $f = 0$ . The approximation space has the form

$$V_N = \{u \in L^2(\Omega) : u|_K \in V_{N,K} \quad \forall K \in \mathcal{T}\}, \quad (51)$$

where the spaces  $V_{N,K}$  are spaces of solutions of the homogeneous Helmholtz equation, e.g., systems of plane waves. For each edge  $e \in \mathcal{E}$ , we select weights  $w_{1,e}$ ,  $w_{2,e} > 0$  and define the functional  $J : V_N \rightarrow \mathbb{R}$  by

$$J(v) := \sum_{e \in \mathcal{E}_I} w_{1,e}^2 \|[v]\|_{L^2(e)}^2 + w_{2,e}^2 \|[\partial_n v]\|_{L^2(e)}^2 + \sum_{e \in \mathcal{E}_B} w_{2,e}^2 \|g - (\partial_n v + \mathbf{i}k v)\|_{L^2(e)}^2;$$

here  $[v]|_e := \llbracket v \rrbracket|_e$  and  $[\partial_n v]|_e := \llbracket \nabla_h v \rrbracket|_e$  represent the jumps of  $v$  and  $\partial_n v$  across the edge  $e$ . If the exact solution  $u$  of (9) is sufficiently regular, then it is a minimizer of  $J$  with  $J(u) = 0$ . In a Least Squares method,  $J$  is minimizer over a finite dimensional space  $V_N$  of the form (51). Its variational form reads:

$$\text{find } u_N \in V_N \text{ s.t. } \langle u_N, v \rangle_{J,N} = \sum_{e \in \mathcal{E}_B} (g, \partial_n v + \mathbf{i}k v)_{L^2(e)} \quad \forall v \in V_N, \quad (52)$$

where

$$\begin{aligned} \langle u, v \rangle_{J,N} := & \sum_{e \in \mathcal{E}_I} w_{1,e}^2 ([u], [v])_{L^2(e)} + w_{2,e}^2 ([\partial_n u], [\partial_n v])_{L^2(e)} + \sum_{e \in \mathcal{E}_B} w_{2,e}^2 (\partial_n u + \mathbf{i}k u, \partial_n v + \mathbf{i}k v)_{L^2(e)}. \end{aligned}$$

The positive semidefinite sesquilinear form  $\langle \cdot, \cdot \rangle_{J,N}$  induces in fact a norm on  $V_N$ : To see the definiteness of  $\langle \cdot, \cdot \rangle_{J,N}$ , we note that  $v \in V_N$  and  $J(v) = 0$  implies that  $v$  is in  $C^1(\Omega)$  and elementwise a solution of the homogeneous Helmholtz equation. Thus, it is a classical solution of the Helmholtz equation on  $\Omega$  and satisfies  $\partial_n v + \mathbf{i}k v = 0$  on  $\partial\Omega$ . The uniqueness assertion for (9) with  $f = 0$  and  $g = 0$  worked out in Example 2.1 then implies  $v = 0$ . Therefore, the minimization problem (52) is well-defined. If the solution  $u$  of (9) satisfies  $u \in H^{3/2+\varepsilon}(\Omega)$  for some  $\varepsilon > 0$ , then  $J(u) = 0$ , and we get quasi-optimality of the Least Squares method in the norm  $\|\cdot\|_{J,N} = J(\cdot)^{1/2}$ :

$$\|u - u_N\|_{J,N}^2 = J(u - u_N) = J(u_N) = \min_{v \in V_N} J(v) = \|u - v\|_{J,N}^2. \quad (53)$$

We mention here that estimates for this minimum can be obtained from (local) estimates in classical Sobolev norms as given in Section 5 using appropriate trace estimates. Turning estimates for  $\|u - u_N\|_{J,N} = J(u_N)^{1/2}$  into estimates in terms of



more familiar norms such as  $\|u - u_N\|_{L^2(\Omega)}$  is not straight forward. It may be expected that the norm of the solution operator of the continuous problem appears again; the next result, which is closely related to [19, 42, 43, 63], illustrates the kind of result one can obtain, in particular in a  $p$ -version setting:

**Lemma 6.2 ([65, Thm. 3.1]).** *Let  $\Omega \subset \mathbb{R}^2$  be a polygon. Let  $w_{1,e} = k$  and  $w_{2,e} = 1$  for all edges and  $g \in L^2(\partial\Omega)$ . Let  $u_N \in V_N$  be the minimizer of  $J$ , where  $V_N$ , given by (51), consists of elementwise solutions of the homogeneous Helmholtz equation.*

(i) *If  $\Omega$  is convex, then  $\|u - u_N\|_{L^2(\Omega)}^2 \leq Ck^{-1} ((kh)^{-1} + (kh)^1) J(u_N)$ .*

(ii) *If  $\Omega$  is not convex, then*

$$\|u - u_N\|_{L^2(\Omega)}^2 \leq Ck^{-1} \left[ (kh)^{-1} + (kh)^1 \left\{ 1 + \min\{1, kh\}^{-2\beta_{\max}} \right\} \right] (C_{\text{sol}}(k) + 1)^2 J(u_N),$$

where  $C_{\text{sol}}(k)$  is defined in (25) and satisfies  $C_{\text{sol}}(k) = O(k^{5/2})$  by Theorem 2.4. The parameter  $\beta_{\max} \geq 0$  can be selected arbitrarily to satisfy the condition  $\beta_{\max} > 1 - \min_i \frac{\pi}{\omega_i}$ , where the  $\omega_i$  are the interior angles of the polygon.

*Proof.* The result (i) is essentially the statement of [65, Thm. 3.1] in a refined form as given in [43, Lemma 3.7]. The statement (ii) is a slightly modification of (i), and we restrict our presentation to that case. The key idea is to obtain  $L^2(\Omega)$ -bounds by a duality argument and use the fact that  $u - u_N$  solves the homogeneous Helmholtz equation elementwise. More precisely, given  $\varphi \in L^2(\Omega)$  we let  $v \in H^1(\Omega)$  solve the adjoint problem

$$-\Delta v - k^2 v = \varphi \quad \text{in } \Omega, \quad \partial_n v - ikv = 0 \quad \text{on } \partial\Omega.$$

By Corollary 3.3, the function  $v$  is in a weighted  $H^2$ -space with

$$\|v\|_{1,k,\Omega} + k^{-1} \|\Phi_{0,\vec{\beta},k} \nabla^2 v\|_{L^2(\Omega)} \leq C(C_{\text{sol}}(k) + 1) \|\varphi\|_{L^2(\Omega)}. \quad (54)$$

Inspection of the arguments underlying the proof of Corollary 3.3 shows that the exponents  $\beta_j \in [0, 1)$  stem from the regularity theory for the Laplacian with Neumann boundary conditions. Hence, in fact  $\beta_j \in [0, 1/2)$  so that  $\nabla_h v$  has an  $L^2$ -trace on all edges of the triangulation (cf. Lemma A.2). For each  $K \in \mathcal{T}$  we then have

$$\|w\|_{L^2(\partial K)}^2 \leq C \left[ h^{-1} \|w\|_{L^2(K)}^2 + h |w|_{H^1(K)}^2 \right] \quad \forall w \in H^1(K), \quad (55)$$

$$\|\nabla w\|_{L^2(\partial K)}^2 \leq C \left[ h^{-1} |w|_{H^1(K)}^2 + h |w|_{H^2(K)}^2 \right] \quad \forall w \in H^2(K), \quad (56)$$

$$\|\nabla w\|_{L^2(\partial K)}^2 \leq C \left[ h^{-1} |w|_{H^1(K)}^2 + h^{1-2\beta} \|r^\beta \nabla^2 w\|_{L^2(K)}^2 \right] \quad \forall w \in H_\beta^{2,2}(K), \quad (57)$$

where, in the last estimate we assume that the origin is at one corner of  $K$  and  $r$  denotes the distance from that corner. These estimates are obtained with the aid of scaling arguments, the multiplicative trace inequality (see [18, Prop. 1.6.3]), and,

in the case of (57) additionally Lemma A.2. From  $ab = \max\{a, b\} \min\{a, b\}$  (for  $a, b \geq 0$ ), we get  $r^\beta = k^{-\beta} (rk)^\beta = k^{-\beta} \min\{1, rk\}^\beta \max\{1, rk\}^\beta$ . Hence, if  $\mathcal{T}_{corner}$  denotes the set of elements that abut on the corners of  $\Omega$ , we get

$$\sum_{K \in \mathcal{T}_{corner}} \|\nabla v\|_{L^2(\partial K)}^2 \leq C \left[ h^{-1} |v|_{H^1(B_h)}^2 + h \min\{1, kh\}^{-2\beta_{max}} \|\Phi_{0, \vec{\beta}, k} \nabla^2 v\|_{L^2(B_h)}^2 \right],$$

where  $B_h = \cup_{K \in \mathcal{T}_{corner}} K$ . Noting that the elements  $K \in \mathcal{T} \setminus \mathcal{T}_{corner}$  are at least  $O(h)$  away from the corners allows us to estimate with (56)

$$\sum_{K \in \mathcal{T} \setminus \mathcal{T}_{corner}} \|\nabla v\|_{L^2(\partial K)}^2 \leq C \left[ h^{-1} |v|_{H^1(\Omega \setminus B_h)}^2 + h^1 \min\{1, kh\}^{-2\beta_{max}} \|\Phi_{0, \vec{\beta}, k} \nabla^2 v\|_{L^2(\Omega \setminus B_h)}^2 \right].$$

Hence, we have the two bounds

$$\begin{aligned} \sum_{K \in \mathcal{T}} \|\nabla v\|_{L^2(\partial K)}^2 &\leq Ch^{-1} |v|_{H^1(\Omega)}^2 + Ch^1 \min\{1, kh\}^{-2\beta_{max}} \|\Phi_{0, \vec{\beta}, k} \nabla^2 v\|_{L^2(\Omega)}^2, \\ \sum_{K \in \mathcal{T}} \|v\|_{L^2(\partial K)}^2 &\leq Ch^{-1} \|v\|_{L^2(\Omega)}^2 + h |v|_{H^1(\Omega)}^2. \end{aligned}$$

Therefore, recalling  $w_{1,e} = k$  and  $w_{2,e} = 1$ , we obtain from these estimates and the *a priori* estimate (54) the bound

$$\begin{aligned} &\sum_{e \in \mathcal{E}} w_{2,e}^{-2} \|v\|_{L^2(e)}^2 + w_{1,e}^{-2} \|\nabla v\|_{L^2(e)}^2 \\ &\leq Ck^{-1} \left[ \frac{1}{kh} + kh \left\{ 1 + \min\{1, kh\}^{-2\beta_{max}} \right\} \right] (C_{sol}(k) + 1)^2 \|\varphi\|_{L^2(\Omega)}^2. \end{aligned} \quad (58)$$

The estimate (58) can be used to bound  $|(u - u_N, \varphi)_{L^2(\Omega)}|$ : Writing the integral as a sum over elements, integrating by parts twice and using that  $u - u_N$  solves the homogeneous Helmholtz equation elementwise yields

$$\begin{aligned} (u - u_N, \varphi)_{L^2(\Omega)} &= \sum_{K \in \mathcal{T}} (u - u_N, -\Delta v - k^2 v)_{L^2(K)} \\ &= \sum_{K \in \mathcal{T}} (\partial_n(u - u_N), v)_{L^2(\partial K)} - \sum_{K \in \mathcal{T}} (u - u_N, \partial_n v)_{L^2(\partial K)} =: \Sigma_1 - \Sigma_2. \end{aligned}$$

The ‘‘DG magic formulas’’ of Lemma 6.1 produce

$$\begin{aligned} \Sigma_1 &= \sum_{e \in \mathcal{E}_I} (\{\{\nabla(u - u_N)\}\}, \llbracket v \rrbracket)_{L^2(e)} + (\llbracket \nabla(u - u_N) \rrbracket, \{\{v\}\})_{L^2(e)} + \sum_{e \in \mathcal{E}_B} (\{\{\nabla(u - u_N)\}\}, \llbracket v \rrbracket)_{L^2(e)} \\ \Sigma_2 &= \sum_{e \in \mathcal{E}_I} (\llbracket u - u_N \rrbracket, \{\{\nabla v\}\})_{L^2(e)} + (\{\{u - u_N\}\}, \llbracket \nabla v \rrbracket)_{L^2(e)} + \sum_{e \in \mathcal{E}_B} (\llbracket u - u_N \rrbracket, \{\{\nabla v\}\})_{L^2(e)}. \end{aligned}$$

For interior edges  $e \in \mathcal{E}_I$  we have  $\llbracket v \rrbracket = 0$  and  $\llbracket \nabla v \rrbracket = 0$  as well as  $\llbracket u \rrbracket = 0$  and  $\llbracket \nabla u \rrbracket = 0$ ; on boundary edges  $e \in \mathcal{E}_B$  we have with the boundary conditions satisfied by  $u$  and  $v$  (i.e.,  $\partial_n u + \mathbf{i}ku = g$  and  $\partial_n v - \mathbf{i}kv = 0$ )

$$(\llbracket \nabla(u - u_N) \rrbracket, \llbracket v \rrbracket)_{L^2(e)} - (\llbracket u - u_N \rrbracket, \llbracket \nabla v \rrbracket)_{L^2(e)} = -((\partial_n + \mathbf{i}k)u_N - g, v)_{L^2(e)}.$$

These observations lead to

$$\begin{aligned} & \left| -(u - u_N, \varphi)_{L^2(\Omega)} \right| = \\ & \left| \sum_{e \in \mathcal{E}_I} (\llbracket u_N \rrbracket, \llbracket \nabla v \rrbracket)_{L^2(e)} + \sum_{e \in \mathcal{E}_I} (\llbracket \nabla u_N \rrbracket, \llbracket v \rrbracket)_{L^2(e)} + \sum_{e \in \mathcal{E}_B} ((\partial_n + \mathbf{i}k)u_N - g, v)_{L^2(e)} \right| \\ & \leq C\sqrt{J(u_N)} \sqrt{\sum_{e \in \mathcal{E}} w_{2,e}^{-2} \|\llbracket v \rrbracket\|_{L^2(e)}^2 + w_{1,e}^{-2} \|\llbracket \nabla v \rrbracket\|_{L^2(e)}^2}, \end{aligned}$$

where, in the last step, we employed the Cauchy-Schwarz inequality for sums. From (58) we therefore get

$$\begin{aligned} & \frac{|(u - u_N, \varphi)_{L^2(\Omega)}|}{\|\varphi\|_{L^2(\Omega)}} \leq \\ & C\sqrt{J(u_N)}k^{-1/2} \left[ (kh)^{-1/2} + (kh)^{1/2} \min\{1, kh\}^{-\beta_{\max}} \right] (C_{\text{sol}}(k) + 1). \end{aligned}$$

Since  $\varphi \in L^2(\Omega)$  is arbitrary, we get the result.  $\square$

**Remark 6.3.** Lemma 6.2 assumes quasi-uniform meshes and the weights  $w_{1,e}$ ,  $w_{2,e}$  do not take the edge length into account. This limits somewhat its applicability in an  $h$ -version context. However, the result is very suitable for a  $p$ -version setting. We point out that in a case where the  $p$ -version features only algebraic rates of convergence, one would have to give the parameters  $w_{1,e}$ ,  $w_{2,e}$  a  $p$ -dependent relative weight as opposed to the situation studied in Lemma 6.2.  $\blacksquare$

### 6.3 Stability of plane wave DG and UWVF

The framework of Discontinuous Galerkin (DG) methods permits another way of deriving numerical schemes that are inherently stable. In a classical, piecewise polynomial setting, this is pursued in [33–35]; related work is in [64]. Here, we concentrate on a setting where the ansatz functions satisfy the homogeneous Helmholtz equation. In particular, we study the plane wave DG method, [36, 43, 63], and the closely related Ultra Weak Variational Formulation (UWVF), [19–21, 46, 55]. We point out that the UWVF can be derived in different way. Here, we follow [19, 36] in viewing it as a special DG method.

Our model problem (9) can be reformulated as a first order system by introducing the flux  $\sigma := (1/\mathbf{i}k)\nabla u$ :

$$\mathbf{i}k\sigma = \nabla u \quad \text{in } \Omega, \quad (59a)$$

$$\mathbf{i}ku - \nabla \cdot \sigma = 0 \quad \text{in } \Omega, \quad (59b)$$

$$\mathbf{i}k\sigma \cdot \mathbf{n} + \mathbf{i}ku = g \quad \text{on } \partial\Omega. \quad (59c)$$

For a mesh  $\mathcal{T}$ , the weak elementwise formulation of (59a), (59b) is for every  $K \in \mathcal{T}$ :

$$\begin{aligned} \int_K \mathbf{i}k \boldsymbol{\sigma} \cdot \bar{\boldsymbol{\tau}} + \int_K u \nabla \cdot \bar{\boldsymbol{\tau}} - \int_{\partial K} u \bar{\boldsymbol{\tau}} \cdot \mathbf{n} &= 0 \quad \forall \boldsymbol{\tau} \in H(\operatorname{div}, K), \\ \int_K \mathbf{i}k u \bar{v} + \int_K \boldsymbol{\sigma} \cdot \nabla \bar{v} - \int_{\partial K} \boldsymbol{\sigma} \cdot \mathbf{n} \bar{v} &= 0 \quad \forall v \in H^1(K), \end{aligned}$$

where  $H(\operatorname{div}, K) = \{u \in L^2(K) : \operatorname{div} u \in L^2(K)\}$  and  $\mathbf{n}$  is the outward pointing normal vector. Replacing the spaces  $H^1(K)$  and  $H(\operatorname{div}, K)$  by finite-dimensional subsets  $V_{N,K} \subset H^1(K)$  and  $\Sigma_{N,K} \subset H(\operatorname{div}, K)$  and, additionally, imposing a coupling between neighboring elements by replacing the multivalued traces  $u$  and  $\boldsymbol{\sigma}$  on the element edges by single-valued numerical fluxes  $\hat{u}_N$ ,  $\hat{\boldsymbol{\sigma}}_N$  to be specified below, leads to the problem: Find  $(u_N, \boldsymbol{\sigma}_N) \in V_{N,K} \times \Sigma_{N,K}$  such that

$$\begin{aligned} \int_K \mathbf{i}k \boldsymbol{\sigma}_N \cdot \bar{\boldsymbol{\tau}} + \int_K u_N \nabla \cdot \bar{\boldsymbol{\tau}} - \int_{\partial K} \hat{u}_N \bar{\boldsymbol{\tau}} \cdot \mathbf{n} &= 0 \quad \forall \boldsymbol{\tau} \in \Sigma_{N,K}, \\ \int_K \mathbf{i}k u_N \bar{v} + \int_K \boldsymbol{\sigma}_N \cdot \nabla \bar{v} - \int_{\partial K} \hat{\boldsymbol{\sigma}}_N \cdot \mathbf{n} \bar{v} &= 0 \quad \forall v \in V_{N,K}. \end{aligned}$$

The variable  $\boldsymbol{\sigma}_N$  can be eliminated by making the assumption that  $\nabla V_{N,K} \subset \Sigma_{N,K}$  for all  $K \in \mathcal{T}$  and then selecting the test function  $\boldsymbol{\tau} = \nabla v$  on each element. This yields after an integration by parts:

$$\int_K \nabla u_N \nabla \bar{v} - k^2 u_N \bar{v} - \int_{\partial K} (u_N - \hat{u}_N) \partial_n \bar{v} - \mathbf{i}k \hat{\boldsymbol{\sigma}}_N \cdot \mathbf{n} \bar{v} = 0 \quad \forall K \in \mathcal{T}. \quad (60)$$

Since  $V_N = \{u \in L^2(\Omega) : u|_K \in V_{N,K} \forall K \in \mathcal{T}\}$  consists of discontinuous functions without any interelement continuity imposed across the element edges, (60) is equivalent to the sum over the elements: Find  $u_N \in V_N$  such that for all  $v \in V_N$

$$\sum_{K \in \mathcal{T}} \int_K \nabla u_N \cdot \nabla \bar{v} - k^2 u_N \bar{v} + \int_{\partial K} (\hat{u}_N - u_N) \nabla \bar{v} \cdot \mathbf{n} - \int_{\partial K} \mathbf{i}k \hat{\boldsymbol{\sigma}}_N \cdot \mathbf{n} \bar{v} = 0. \quad (61)$$

This formulation is now completed by specifying the fluxes  $\hat{u}_N$  and  $\hat{\boldsymbol{\sigma}}_N$ , which at the same time takes care of the boundary condition (59c):

- For interior edges  $e \in \mathcal{E}_I$

$$\hat{\boldsymbol{\sigma}}_N = \frac{1}{\mathbf{i}k} \{\{\nabla_h u_N\}\} - \alpha \llbracket u_N \rrbracket, \quad \hat{u}_N = \{\{u_N\}\} - \beta \frac{1}{\mathbf{i}k} \llbracket \nabla_h u_N \rrbracket. \quad (62a)$$

- For boundary edges  $e \in \mathcal{E}_B$

$$\hat{\boldsymbol{\sigma}}_N = \frac{1}{\mathbf{i}k} \nabla_h u_N - \frac{1-\delta}{\mathbf{i}k} (\nabla_h u_N + \mathbf{i}k u_N \mathbf{n} - g \mathbf{n}). \quad (62b)$$

$$\hat{u}_N = u_N - \frac{\delta}{\mathbf{i}k} (\nabla_h u \cdot \mathbf{n} + \mathbf{i}k u_N - g). \quad (62c)$$

Different choices of the parameters  $\alpha$ ,  $\beta$ ,  $\delta$  lead to different methods analyzed in the literature. For example:

1.  $\alpha = \beta = \delta = 1/2$ : this is the UWVF as analyzed in [19–21, 46, 55] if the spaces  $V_{N,K}$  consist of a space  $W_{PW}^p$  of plane waves.
2.  $\alpha = O(p/(kh \log p))$ ,  $\beta = O((kh \log p)/p)$ ,  $\delta = O((kh \log p)/p)$ : this choice is introduced and advocated in [43, 63] in conjunction with  $V_{N,K} = W_{PW}^p$ .

With these choices of fluxes, the formulation (61) takes the form

$$\text{Find } u_N \in V_N \text{ s.t. } A_N(u_N, v) = l(v) \quad \forall v \in V_N, \quad (63)$$

where the sesquilinear form  $A_N$  and the linear form  $l$  are given by

$$\begin{aligned} A_N(u, v) &= \int_{\Omega} \nabla_h u \cdot \nabla_h \bar{v} - k^2 u \bar{v} - \int_{\mathcal{E}_I} [u] \{ \nabla_h \bar{v} \} - \int_{\mathcal{E}_I} \{ \nabla_h u \} [\bar{v}] - \int_{\mathcal{E}_B} \delta u \partial_n \bar{v} - \int_{\mathcal{E}_B} \delta \partial_n u \bar{v} \\ &\quad - \frac{1}{ik} \int_{\mathcal{E}_I} \beta [\nabla_h u] [\nabla_h \bar{v}] - \frac{1}{ik} \int_{\mathcal{E}_B} \delta \partial_n u \partial_n \bar{v} + ik \int_{\mathcal{E}_I} \alpha [u] [\bar{v}] + ik \int_{\mathcal{E}_B} (1 - \delta) u \bar{v} \quad (64) \\ l(v) &= -\frac{1}{ik} \int_{\mathcal{E}_B} \delta g \partial_n \bar{v} + \int_{\mathcal{E}_B} (1 - \delta) g \bar{v}. \end{aligned}$$

So far, the choice of the spaces  $V_{N,K}$  is arbitrary. If the approximation spaces  $V_{N,K}$  (more precisely: the test spaces) consist of piecewise solutions of the homogeneous Helmholtz equation, then a further integration by parts is possible to eliminate all volume contributions in  $A_N$ . Indeed, Lemma 6.1 produces

$$\sum_{K \in \mathcal{T}} \int_K \nabla u \cdot \nabla \bar{v} - k^2 u \bar{v} = \sum_{K \in \mathcal{T}} \int_{\partial K} u \nabla \bar{v} \mathbf{n} = \int_{\mathcal{E}_I} [u] \{ \nabla \bar{v} \} + \{ u \} [\nabla \bar{v}] + \int_{\mathcal{E}_B} [u] \{ \nabla \bar{v} \}$$

so that  $A_N$  simplifies to

$$\begin{aligned} A_N(u, v) &= \int_{\mathcal{E}_I} \{ u \} [\nabla_h \bar{v}] + \frac{1}{k} \int_{\mathcal{E}_I} \beta [\nabla_h u] [\nabla_h \bar{v}] - \int_{\mathcal{E}_I} \{ \nabla_h u \} [\bar{v}] + ik \int_{\mathcal{E}_I} \alpha [u] [\bar{v}] \\ &\quad + \int_{\mathcal{E}_B} (1 - \delta) u \partial_n \bar{v} + \frac{1}{k} \int_{\mathcal{E}_B} \delta \partial_n u \partial_n \bar{v} - \int_{\mathcal{E}_B} \delta \partial_n u \bar{v} + ik \int_{\mathcal{E}_B} (1 - \delta) u \bar{v}. \end{aligned}$$

Next, we make the important observation that  $\text{Im} A_N$  induces a norm on the space  $V_N$  if  $\alpha, \beta > 0$  and  $\delta \in (0, 1)$ . Indeed:

1.  $\alpha, \beta > 0$  and  $\delta \in (0, 1)$  implies  $\text{Im} A_N(v, v) \geq 0 \quad \forall v \in V_N$  by inspection of (64).
2.  $\text{Im} A_N(v, v) = 0$  and the fact that  $V_N$  consists of elementwise solutions of the homogeneous Helmholtz equation implies as in the case of  $\langle \cdot, \cdot \rangle_{J,N}$  in Section 6.2 that  $v \in C^1(\Omega)$  solves the homogeneous Helmholtz equation and  $\partial_n v = v = 0$  on  $\partial\Omega$ ; the uniqueness assertion of Example 2.1 then proves  $v \equiv 0$ .

This is at the basis of the convergence analysis. Introducing

$$\begin{aligned}
\|u\|_{DG}^2 &:= \sqrt{\operatorname{Im} A_N(u, u)} = \frac{1}{k} \|\beta^{1/2} \llbracket \nabla_h u \rrbracket\|_{L^2(\mathcal{E}_I)}^2 + \|\alpha^{1/2} \llbracket u \rrbracket\|_{L^2(\mathcal{E}_I)}^2 \\
&\quad + \frac{1}{k} \|\delta^{1/2} \partial_n u\|_{L^2(\mathcal{E}_B)}^2 + k \|(1 - \delta)^{1/2} u\|_{L^2(\mathcal{E}_B)}^2, \\
\|u\|_{DG,+}^2 &:= \|u\|_{DG}^2 + k \|\beta^{-1/2} \llbracket u \rrbracket\|_{L^2(\mathcal{E}_I)}^2 + k^{-1} \|\alpha^{-1/2} \llbracket u \rrbracket\|_{L^2(\mathcal{E}_I)}^2 + k \|\delta^{-1/2} u\|_{L^2(\mathcal{E}_B)}^2,
\end{aligned}$$

we can formulate coercivity and continuity results:

**Proposition 6.4** ([19, 43]). *Let  $V_N$  consist of piecewise solutions of the homogeneous Helmholtz equation. Then  $\|\cdot\|_{DG}$  is a norm on  $V_N$  and for some  $C > 0$  depending solely on the choice of  $\alpha$ ,  $\beta > 0$ , and  $\delta \in (0, 1)$ :*

$$\begin{aligned}
\operatorname{Im} A_N(u, u) &= \|u\|_{DG}^2 \quad \forall u \in V_N, \\
|A_N(u, v)| &\leq C \|u\|_{DG,+} \|v\|_{DG} \quad \forall u, v \in V_N
\end{aligned}$$

Let the solution of  $u$  of (9) (with  $f = 0$ ) satisfy  $u \in H^{3/2+\varepsilon}(\Omega)$  for some  $\varepsilon > 0$ . Then, by consistency of  $A_N$ , the solution  $u_N \in V_N$  of (63) satisfies the following quasi-optimality estimate for some  $C > 0$  independent of  $k$ :

$$\|u - u_N\|_{DG} \leq C \inf_{v \in V_N} \|u - v\|_{DG,+}. \quad (65)$$

Several comments are in order:

1. The UWVF of [20] featured quasi-optimality in a residual type norm. We recall that the UWVF is a DG method for the particular choice  $\alpha = \beta = \delta = 1/2$ .
2. When  $V_N$  consists (elementwise) of systems of plane waves or generalized harmonic polynomials, then the infimum in (65) can be estimated using approximation results on the elements by taking appropriate traces. This is worked out in detail in [42, 43, 63] and earlier in an  $h$ -version setting in [20] (see also [19]).
3. The  $\|\cdot\|_{DG}$ -norm controls the error on the skeleton  $\mathcal{E}$  only. The proof of Lemma 6.2 shows how error estimates in such norms can be used to obtain estimates for  $\|u - u_N\|_{L^2(\Omega)}$ ; we refer again to [19] where this worked out for the UWVF and to [42, 43, 63] where the case of the plane wave DG is studied. As pointed out in Remark 6.3, quasi-uniformity of the underlying mesh  $\mathcal{T}$  is an important ingredient for the arguments of Lemma 6.2.

It is noteworthy that Proposition 6.4 does not make any assumptions on the mesh size  $h$  and the space  $V_N$  except that it consist of piecewise solutions of the homogeneous Helmholtz equation. Optimal error estimates are possible in an  $h$ -version setting, where the number of plane waves per element is kept fixed:

**Proposition 6.5** ([36]). *Let  $\Omega$  be convex. Assume that  $V_{N,K} = W_{P_w}^{2\mu+1}$  ( $\mu \geq 1$  fixed) for all  $K \in \mathcal{T}$ . Assume that  $\alpha$  is of the form  $\alpha = \mathbf{a}/(kh)$  and that  $\beta > 0$ ,  $\delta \in (0, 1/2)$ . Then there exist  $\mathbf{a}_0$ ,  $c_0$ ,  $C > 0$  (all independent of  $h$  and  $k$ ) such that if  $\mathbf{a} \geq \mathbf{a}_0$  and  $k^2 h \leq c_0$ , then following error bound is true:*

$$\|u - u_N\|_{1,DG} \leq C \inf_{v \in V_N} \|u - v\|_{1,DG,+};$$

here,  $\|\cdot\|_{1,DG}$  and  $\|\cdot\|_{1,DG,+}$  are given by  $\|v\|_{1,DG}^2 := \sum_{K \in \mathcal{T}} |v|_{H^1(K)}^2 + k^2 \|v\|_{L^2(K)}^2 + \|v\|_{DG}^2$  and  $\|v\|_{1,DG,+}^2 := \sum_{K \in \mathcal{T}} |v|_{H^1(K)}^2 + k^2 \|v\|_{L^2(K)}^2 + \|v\|_{DG,+}^2$ .

*Proof.* The proof follows by inspection of the procedure in [36, Sec. 5] and is stated in [63, Props. 4.2, 4.3]. The essential ingredients of the proof are: (a) inverse estimates for systems of plane waves that have been made in available in [36] so that techniques of standard DG methods can be used to treat  $A_N$ ; (b) use of duality arguments as in Lemma 4.1 to treat the  $L^2$ -norm of the error; (c) the fact that in an  $h$ -version setting, plane waves have some approximation power for arbitrary functions in  $H^2$  (this is analogous to Lemma 4.8).  $\square$

## 7 Remarks on 1D

The 1D situation is rather special in that pollution can be completely eliminated; the underlying reason is that the space of solutions of the homogeneous Helmholtz equation is finite-dimensional (two dimensional, in fact). We illustrate this for the following model problem:

$$-u'' - k^2 u = f \quad \text{in } \Omega = (0, 1), \quad u(0) = 0, \quad u'(1) - iku(1) = 0. \quad (66)$$

Let  $\mathcal{T}$  be a mesh on  $\Omega$  with nodes  $0 = x_0 < x_1 < \dots < x_N = 1$ . We assume that the mesh size  $h := \max_i (x_{i+1} - x_i)$  satisfies  $kh < \pi$ . For each node  $x_i$ , let  $\psi_i \in H^1(\Omega)$  be defined by the conditions

$$\psi_i(x_j) = \delta_{ij}, \quad (-\psi_i'' - k^2 \psi_i)|_K = 0 \quad \forall K \in \mathcal{T}$$

and let  $V_N^{opt} = \text{span}\{\psi_i : i = 1, \dots, N\}$ . Thus, the functions  $\psi_i$  are piecewise solutions of the homogeneous Helmholtz equation. The Galerkin method based on  $V_N^{opt}$  is:

$$\text{Find } u_N \in V_N^{opt} \text{ s.t. } \int_{\Omega} u_N' \bar{v}' - k^2 u_N \bar{v} - iku_N(1) \bar{v}(1) = \int_{\Omega} f \bar{v} \quad \forall v \in V_N^{opt}. \quad (67)$$

The Galerkin method based on  $V_N^{opt}$  is nodally exact:

**Lemma 7.1.** *There exist constants  $C_1, C_2 > 0$  independent of  $k$  such that the following is true for  $kh \leq C_1$ :*

- (i) *The functions  $\psi_i$  are well-defined.*
- (ii) *The method (67) is nodally exact.*
- (iii) *For  $f \in L^2(\Omega)$  there holds  $\|u - u_N\|_{1,k,\Omega} \leq C_2(hk) \|f\|_{L^2(\Omega)}$ .*

*Proof.* Elementary considerations show that for  $kh < \pi$ , the functions  $\psi_i$  are well-defined.

The most interesting feature of Lemma 7.1 is the nodal exactness. To that end, we note that the Green's function for (66) is

$$G(x, y) = \frac{1}{k} \begin{cases} \sin kx e^{\mathbf{i}ky} & 0 < x < y \\ \sin ky e^{\mathbf{i}kx} & y < x < 1. \end{cases}$$

Let  $e \in H^1(\Omega)$  satisfy  $e(0) = 0$  and the Galerkin orthogonality condition

$$C(e, v) := \int_{\Omega} e' \bar{v}' - k^2 e \bar{v} - \mathbf{i}ke(1) \bar{v}(1) = 0 \quad \forall v \in V_N^{opt}. \quad (68)$$

The key observation is that for each  $x_i, i = 1, \dots, N$ , the function  $v_i := G(\cdot, x_i) \in V_N^{opt}$  since is a solution of the homogeneous Helmholtz equation on  $(0, x_i) \cup (x_i, 1)$ , it satisfies  $G(0, x_i) = 0$ . Furthermore, we have  $v_i'(x_N) - \mathbf{i}kv_i(x_N) = 0$ . Hence, we get from the Galerkin orthogonality (68) by an integration by parts:

$$\begin{aligned} 0 &= \int_{\Omega} e' \bar{v}_i' - k^2 e \bar{v}_i + \mathbf{i}ke(1) \bar{v}_i(1) \\ &= \int_0^{x_i} e' \bar{v}_i' - k^2 e \bar{v}_i + \int_{x_i}^{x_N} e' \bar{v}_i' - k^2 e \bar{v}_i - \mathbf{i}ke(1) \bar{v}_i(1) \\ &= \int_0^{x_i} e(-\bar{v}_i'' - k^2 \bar{v}_i) + e(x_i) \bar{v}_i'(x_i) - e(x_i) \bar{v}_i'(x_i) \\ &\quad + \int_{x_i}^1 e(-\bar{v}_i'' - k^2 \bar{v}_i) + e(1) \bar{v}_i'(1) - \mathbf{i}ke(1) \bar{v}_i(1) \\ &= e(x_i) [\bar{v}_i'](x_i) + e(1) \bar{v}_i'(1) - \mathbf{i}ke(1) \bar{v}_i(1) \\ &= e(x_i) [\bar{v}_i'](x_i); \end{aligned}$$

here, we have employed the standard notation for the jump of a piecewise smooth function  $w$ :  $[w'](x_i) := \lim_{x \rightarrow x_i^-} w'(x) - \lim_{x \rightarrow x_i^+} w'(x)$ . Since  $[v_i'](x_i) \neq 0$ , we conclude

$$e(x_i) = 0 \quad \forall i \in \{1, \dots, N\}.$$

Hence, the FEM (67) is nodally exact.

The above argument also shows that any  $e \in V_N^{opt}$  satisfying (68) must satisfy  $e(x_i) = 0$  for all  $i \in \{1, \dots, N\}$ . Hence, by the definition of  $V_N^{opt}$  as the span of the functions  $\psi_i$ , we conclude  $e \equiv 0$ . Thus, the kernel of the linear systems described by (67) is trivial. By the usual dimension argument, we have unique solvability.

We have the *a priori* bound

$$\|u\|_{1,k,\Omega} \leq C \|f\|_{L^2(\Omega)} \quad (69)$$

for the solution  $u$  of (66) (as for the model problem (9), this can be shown using the test function  $v = xu'$ ; an alternative proof based on the Green's function and the representation

$$u(x) = \int_{\Omega} G(x, y) f(y) dy$$

is given in [47, Thm. 4.4]). If one denotes by  $\varphi_i$  the classical piecewise linear hat function associated with node  $x_i$ , then one has by Taylor expansion



$$\|\varphi_i - \psi_i\|_{L^\infty(K)} \leq C(kh_K)^2, \quad \|(\varphi_i - \psi_i)'\|_{L^\infty(K)} \leq Ck^2h_K, \quad \forall K \in \mathcal{T}, \quad (70)$$

where  $h_K = \text{diam}K$ . The approximation properties follow easily from the nodal exactness. Specifically, denoting by  $Iu$  the classical piecewise linear interpolant of  $u$  and by  $\tilde{I}u \in V_N^{opt}$  the nodal interpolant determined by  $\tilde{I}u(x_i) = u(x_i)$  for all  $i \in \{0, \dots, N\}$ , we have the well-known estimate

$$\|u - Iu\|_{1,k,\Omega} \leq C(kh^2 + h)\|u''\|_{L^2(\Omega)} \leq Ckh(1 + kh)\|f\|_{L^2(\Omega)},$$

where we used the differential equation and the bound (69) to estimate  $\|u''\|_{L^2(\Omega)} \leq \|f\|_{L^2(\Omega)} + k^2\|u\|_{L^2(\Omega)} \leq Ck\|f\|_{L^2(\Omega)}$ . Next, we estimate the difference  $Iu - \tilde{I}u$ . The multiplicative trace inequality takes the form

$$h_K\|w\|_{L^\infty(K)}^2 \leq C\left[\|w\|_{L^2(K)}^2 + h_K\|w\|_{L^2(K)}\|w'\|_{L^2(K)}\right] \quad \forall w \in H^1(K). \quad (71)$$

Hence, the estimates (70), (71) imply

$$\begin{aligned} \|Iu - \tilde{I}u\|_{1,k,\Omega}^2 &\leq C \sum_{K \in \mathcal{T}} k^2(kh_K)^2(1 + (kh_K)^2)h_K\|u\|_{L^\infty(K)}^2 \\ &\leq C(kh)^2(1 + (kh)^2) \left[ k^2\|u\|_{L^2(\Omega)}^2 + k^2h^2\|u'\|_{L^2(\Omega)}^2 \right] \\ &\leq C(kh)^2(1 + (kh)^2)^2\|u\|_{1,k,\Omega}^2. \end{aligned}$$

An appeal to (69) concludes the argument.  $\square$

Several comments are in order concerning the stability of the method:

1. In the 1D situation, the good stability properties of high order Galerkin FEM can alternatively be understood in light of Lemma 7.1: Applying the Galerkin method to a classical high order method and then condensing out the degrees of freedom corresponding to internal shape functions (“bubbles”), leads to a linear system that is identical to the one obtained by using shape functions  $\psi_i^p$ ,  $i = 0, \dots, N$ , that satisfy  $\psi_i^p(x_j) = \delta_{ij}$  and additionally

$$\int_K (\psi_i^p)' \bar{v}' - k^2 \psi_i^p \bar{v} = 0 \quad \forall v \in H_0^1(K) \cap \mathcal{P}_p$$

Since on a fixed mesh  $\mathcal{T}$ , we have  $\lim_{p \rightarrow \infty} \psi_i^p = \psi_i$ , better stability properties of higher order methods may reasonably be hoped for.

2. The system matrix of the Galerkin FEM based on the space  $V_N^{opt}$  is a tridiagonal matrix. The same matrix can also be obtained in different ways. Consider, for example, the sesquilinear form

$$B(u, v) := \int_\Omega u' \bar{v}' - ku \bar{v} - iku(1) \bar{v}(1) + \sum_{K \in \mathcal{T}} \tau_K \int_\Omega L_k u L_k \bar{v}, \quad (72)$$

where  $L_k = -\frac{d^2}{dx^2} - k^2$ . For a suitable choice of the parameters  $\tau_K$  in dependence on  $k$  and  $h_K$ , the system matrix resulting from this  $B$  using the classical piecewise linear hat functions leads to the same matrix as the Galerkin method based on the shape functions  $\psi_i$ ,  $i = 1, \dots, N$ . In 1D, it is therefore possible to design nodally exact methods based on the stabilization techniques in the form (72). In [12], a nodally exact method is derived using other techniques.

**Acknowledgement:** Financial support by the *Vienna Science and Technology Fund* (WWTF) is gratefully acknowledged.

## Appendix

For the reference triangle  $\widehat{K} := \{(x, y) : 0 < x < 1, 0 < y < 1 - x\}$  and  $\beta \in [0, 1)$  the following two lemmas require the spaces  $H_\beta^{1,1}(\widehat{K})$ ,  $H_\beta^{2,2}(\widehat{K})$  as well as the Besov spaces  $B_{2,\infty}^s(\widehat{K})$ . The spaces  $B_{2,\infty}^s(\widehat{K})$  are defined by interpolation using the  $K$ -functional (see, e.g., [18, Chap. 12]). For  $m \in \{1, 2\}$ , the spaces  $H_\beta^{m,m}(\widehat{K})$  are determined by the norm  $\|u\|_{H_\beta^{m,m}(\widehat{K})}^2 := \|u\|_{H^{m-1}(\widehat{K})}^2 + \|r^\beta \nabla^m u\|_{L^2(\widehat{K})}^2$ , where  $r$  denotes the distance from the origin.

**Lemma A.1.** *Let  $\widehat{K}$  be the reference triangle. Let  $\beta \in [0, 1)$ . Then the embeddings  $H_\beta^{2,2}(\widehat{K}) \subset B_{2,\infty}^{2-\beta}(\widehat{K})$  and  $H_\beta^{1,1}(\widehat{K}) \subset B_{2,\infty}^{1-\beta}(\widehat{K})$  are continuous. The embeddings  $H_\beta^{2,2}(\widehat{K}) \subset H^{2-\beta-\varepsilon}(\widehat{K})$  and  $H_\beta^{1,1}(\widehat{K}) \subset H^{1-\beta-\varepsilon}(\widehat{K})$  are compact for all  $\varepsilon > 0$ .*

*Proof.* Since the case  $\beta = 0$  corresponds to classical Sobolev spaces, we restrict our attention here to the situation  $\beta \in (0, 1)$ . The argument follows ideas presented in [11, Thm. 2.1] and [10]. We start with the following two Hardy inequalities for sufficiently smooth functions

$$\|r^{\beta-1} \nabla u\|_{L^2(\widehat{K})} \leq C \|u\|_{H_\beta^{2,2}(\widehat{K})}, \quad (\text{A.1})$$

$$\|r^{\beta-2} (u - u(0))\|_{L^2(\widehat{K})} \leq C \|u\|_{H_\beta^{2,2}(\widehat{K})}; \quad (\text{A.2})$$

here, (A.1) is shown, for example, in [57, Lemma A.1.7] and (A.2) follows from combining [10, Lemma 4.2] with (A.1). Noting that [10, (2.2)] states the continuous embedding  $H_\beta^{2,2}(\widehat{K}) \subset C(\widehat{K})$ , we have that  $u(0)$  in (A.2) is indeed well-defined.

We employ the real method of interpolation and write  $B_{2,\infty}^{2-\beta} = (L^2, H^2)_{1-\beta/2, \infty}$ . Our method of proof consists in showing that for  $\theta = 1 - \beta/2$  we have

$$\sup_{t \in (0,1)} t^{-\theta} K(t, \tilde{u}) \leq C \|u\|_{H_\beta^{2,2}(\widehat{K})}, \quad \tilde{u} := u - u(0),$$

for some  $C > 0$  independent of  $u$ . To that end, we proceed as in the proof of [10, Lemma 2.1]. For every  $\delta > 0$ , let  $\chi_\delta \in C_0^\infty(\mathbb{R}^2)$  with  $\chi \equiv 1$  on  $B_{\delta/2}(0)$  and  $\text{supp } \chi \subset$

$\chi_\delta B_\delta(0)$  as well as  $\|\nabla^j \chi_\delta\|_{L^\infty(\mathbb{R}^2)} \leq C\delta^{-j}$ ,  $j \in \{0, 1, 2\}$ . We define the splitting

$$\tilde{u} = \chi_\delta \tilde{u} + (1 - \chi_\delta) \tilde{u} =: u_1 + u_2$$

Then from (A.1) and (A.2)

$$\begin{aligned} \|\chi_\delta \tilde{u}\|_{L^2(\hat{K})} &\leq C \|\tilde{u}\|_{L^2(\hat{K} \cap B_\delta(0))} \leq \delta^{2-\beta} \|r^{\beta-2} \tilde{u}\|_{L^2(\hat{K})} \leq C \delta^{2-\beta} \|u\|_{H_\beta^{2,2}(\hat{K})}, \\ |(1 - \chi_\delta) \tilde{u}|_{H^2(\hat{K})} &\leq \\ C \delta^{-2} \|\tilde{u}\|_{L^2((\hat{K} \cap B_\delta(0)) \setminus B_{\delta/2}(0))} &+ C \delta^{-1} \|\nabla \tilde{u}\|_{L^2((\hat{K} \cap B_\delta(0)) \setminus B_{\delta/2}(0))} + C \|\nabla^2 \tilde{u}\|_{L^2(\hat{K} \setminus B_{\delta/2}(0))} \\ &\leq C \delta^{-2+2-\beta} \|r^{\beta-2} \tilde{u}\|_{L^2(\hat{K})} + C \delta^{-1+1-\beta} \|r^{\beta-1} \nabla \tilde{u}\|_{L^2(\hat{K})} + C \delta^{-\beta} \|r^\beta \nabla^2 \tilde{u}\|_{L^2(\hat{K})} \\ &\leq C \delta^{-\beta} \|u\|_{H_\beta^{2,2}(\hat{K})}. \end{aligned}$$

From this, we can infer for any  $\delta \in (0, 1)$

$$K(t, \tilde{u}) \leq \|u_1\|_{L^2(\hat{K})} + t \|u_2\|_{H^2(\hat{K})} \leq C \|u\|_{H_\beta^{2,2}(\hat{K})} \left[ \delta^{2-\beta} + t \delta^{-\beta} \right].$$

Selecting  $\delta = t^{1/2}$  gives  $K(t, \tilde{u}) \leq C t^{1-\beta/2} \|u\|_{H_\beta^{2,2}(\hat{K})}$ . Finally, the compactness assertions of the embeddings follows from the compactness of the embeddings  $B_{2,\infty}^s \subset B_{2,2}^{s'} = H^{s'}$  for  $s' < s$ .  $\square$

**Lemma A.2.** *Let  $\beta \in [0, 1/2)$  and  $\hat{K}$  be the reference triangle. Then there exists  $C > 0$  such that for all  $u \in H_\beta^{1,1}(\hat{K})$  there holds  $\|u\|_{L^2(\partial\hat{K})} \leq C \left[ \|u\|_{L^2(\hat{K})} + \|r^\beta \nabla u\|_{L^2(\hat{K})} \right]$ .*

*Proof.* For each  $s > 1/2$ , we have the inequality  $\|u\|_{L^2(\partial\hat{K})} \leq C_s \|u\|_{H^s(\hat{K})}$ . From the embedding  $H_\beta^{1,1}(\hat{K}) \subset H^{1-\beta}(\hat{K})$  of Lemma A.1, we then get  $\|u\|_{L^2(\partial\hat{K})} \leq C \|u\|_{H^s(\hat{K})} \leq C \left[ \|u\|_{L^2(\hat{K})} + \|r^\beta \nabla u\|_{L^2(\hat{K})} \right]$ .  $\square$

**Lemma A.3.** *Let  $\beta \in [0, 1)$  and  $\Omega \subset \mathbb{R}^2$  be a finite sector with apex at the origin. Let  $u \in C^\infty(\Omega)$  satisfy*

$$\|\Phi_{n,\beta,1} \nabla^{n+2} u\|_{L^2(\Omega)} \leq C_u \gamma_u n! \quad \forall n \in \mathbb{N}_0.$$

*Then, for  $k \geq k_0 > 0$ , there exist constants  $C, \gamma > 0$  (depending only on  $\beta, \Omega, \gamma_u$  and  $k_0$ ) such that*

$$\|\Phi_{n,\beta,k} \nabla^{n+2} u\|_{L^2(\Omega)} \leq C C_u k^{-(2-\beta)} \gamma^n \max\{n, k\}^{n+2} \quad \forall n \in \mathbb{N}_0.$$

*Proof.* Lemma A.4 yields

$$\frac{1}{\max\{n, k\}^{n+2}} \Phi_{n,\beta,k}(x) \leq C k^{-(2-\beta)} \gamma^n \frac{1}{n!} \Phi_{n,\beta,1}(x) \quad \forall x \in \Omega,$$

where  $C, \gamma > 0$  are independent of  $n$  and  $k$ . The result now follows.  $\square$

**Lemma A.4.** *Let  $\beta \in [0, 1)$ . Then for  $0 < r < R$  and all  $n \in \mathbb{N}_0$*

$$\min \left( 1, \frac{r}{\min \left\{ 1, \frac{n+1}{k+1} \right\}} \right)^{n+\beta} \frac{1}{\max \{n, k\}^{n+2}} \leq C k^{-(2-\beta)} \gamma^n r^{n+\beta} \frac{1}{n^{n+2}}$$

*Proof.* We denote the left-hand side by  $lhs$  and consider several cases.

*case 1:*  $n \leq k$  and  $r(k+1) \leq n+1$ :

$$\begin{aligned} lhs &= \left( \frac{(k+1)r}{n+1} \right)^{n+\beta} \frac{1}{k^{n+2}} \\ &= r^{n+\beta} \frac{1}{n^{n+2}} \left( \frac{n}{n+1} \right)^{n+2} (n+1)^{2-\beta} \left( \frac{k+1}{k} \right)^{n+2} (k+1)^{-(2-\beta)} \\ &\leq C \gamma^n k^{-(2-\beta)} r^{n+\beta} \frac{1}{n^{n+2}} \end{aligned}$$

for suitable  $C, \gamma > 0$  if we assume that  $k \geq k_0 > 0$ .

*case 2:*  $n \leq k$  and  $r(k+1) > n+1$ :

$$\begin{aligned} lhs &= \frac{1}{k^{n+2}} = \frac{1}{k^{2-\beta}} \frac{1}{k^{n+\beta}} = \frac{1}{k^{2-\beta}} \left( \frac{k+1}{k} \right)^{n+\beta} \frac{1}{(k+1)^{n+\beta}} \\ &\leq \frac{1}{k^{2-\beta}} \left( \frac{k+1}{k} \right)^{n+\beta} \left( \frac{r}{n+1} \right)^{n+\beta} \\ &\leq C \gamma^n k^{-(2-\beta)} r^{n+\beta} \frac{1}{n^{n+2}} \end{aligned}$$

for suitable  $C, \gamma > 0$ .

*case 3:*  $n > k$ : Then, for  $0 < r < R$

$$\begin{aligned} lhs &= (\min \{1, r\})^{n+\beta} \frac{1}{n^{n+2}} \leq r^{n+\beta} \frac{1}{n^{n+2}} \leq k^{-(2-\beta)} r^{n+\beta} \frac{1}{n^{n+2}} n^{2-\beta} \\ &\leq C k^{-(2-\beta)} r^{n+\beta} \frac{1}{n^{n+2}} \gamma^n \end{aligned}$$

for suitable  $C, \gamma > 0$ .  $\square$

## Further results and proofs

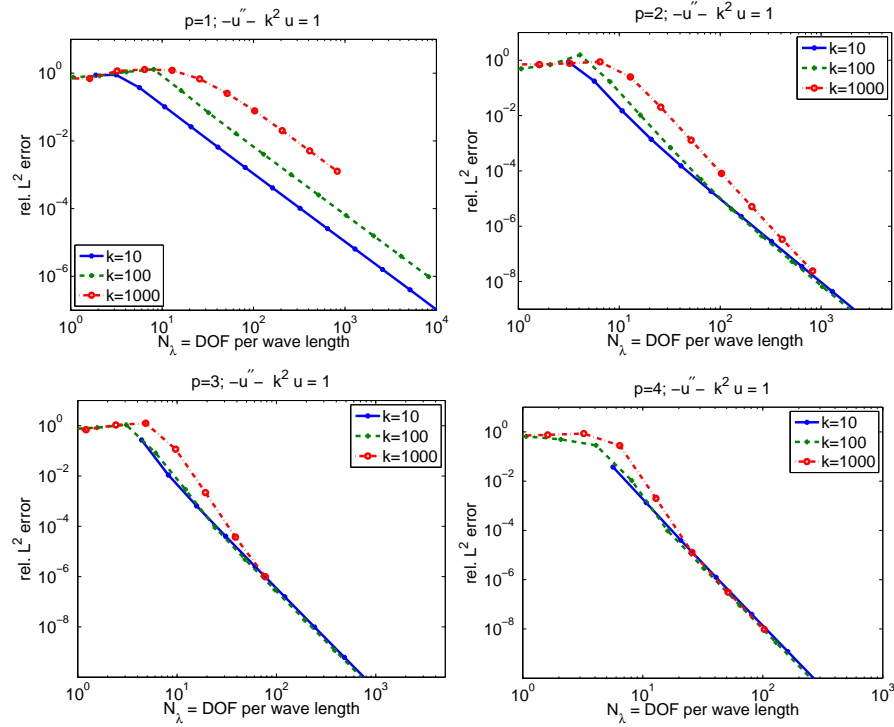
**Example A.5.** In Example 1.1, we studied the convergence behavior of the  $h$ -FEM in the  $H^1(\Omega)$ -seminorm. In Fig. 4 we present the corresponding results for the convergence in the  $L^2(\Omega)$ -norm by plotting  $\|u - u_N\|_{L^2(\Omega)} / \|u\|_{L^2(\Omega)}$  vs. the number of degrees of freedom per wavelength  $N_\lambda$ . For  $p = 1$ , we observe

$$\frac{\|u - u_N\|_{L^2}}{\|u\|_{L^2}} \approx CkN_\lambda^{-2}, \quad N_\lambda \rightarrow \infty,$$

which is in agreement with the analysis given in [47, Sec. 4.6.4]. The cases  $p > 1$  seem to behave differently as we observe

$$\frac{\|u - u_N\|_{L^2}}{\|u\|_{L^2}} \approx CN_\lambda^{-(p+1)}, \quad N_\lambda \rightarrow \infty$$

■



**Fig. 4** Performance of  $h$ -FEM for (2). Top:  $p = 1, p = 2$ . Bottom:  $p = 3, p = 4$  (cf. Examples A.5, 1.1).

*Proof of Remark 4.9.* By interpolation using the  $K$ -functional we can write  $H^{1+\theta} = (H^1, H^2)_{\theta,2}$  for  $\theta \in (0, 1)$ . Hence, every  $u \in H^{1+\theta}(\Omega)$  can be decomposed as  $u = u_1 + u_2$  with

$$\|u_1\|_{H^1(\Omega)} \leq t^\theta \|u\|_{H^{1+\theta}(\Omega)}, \quad \|u_2\|_{H^2(\Omega)} \leq t^{\theta-1} \|u\|_{H^{1+\theta}(\Omega)}, \quad (\text{A.3})$$

where  $t > 0$  is arbitrary. The proof of Lemma 4.8 shows that  $u_1$  and  $u_2$  can be approximated from  $V_N$  as follows:

$$\begin{aligned} \inf_{v \in V_N} \|u_2 - v\|_{L^2(\Omega)} + h \|\nabla(u_2 - v)\|_{L^2(\Omega)} &\leq h^2 \|u_2\|_{H^2(\Omega)} + (kh)^2 \|u_2\|_{L^2(\Omega)} \\ \inf_{v \in V_N} \|u_1 - v\|_{L^2(\Omega)} + h \|\nabla(u_1 - v)\|_{L^2(\Omega)} &\leq h \|u_1\|_{H^1(\Omega)} + (kh)^2 \|u_1\|_{L^2(\Omega)}. \end{aligned}$$

Using  $t = h$  in (A.3) we therefore get

$$\inf_{v \in V_N} \|u - v\|_{L^2(\Omega)} + h \|\nabla(u - v)\|_{L^2(\Omega)} \leq h^{1+\theta} \|u\|_{H^{1+\theta}(\Omega)} + (kh)^2 \left[ \|u_1\|_{L^2(\Omega)} + \|u_2\|_{L^2(\Omega)} \right].$$

The decomposition  $u = u_1 + u_2$  and the triangle inequality yield  $\|u_1\|_{L^2(\Omega)} + \|u_2\|_{L^2(\Omega)} \leq \|u\|_{L^2(\Omega)} + 2\|u_1\|_{L^2(\Omega)} \leq \|u\|_{L^2(\Omega)} + 2\|u_1\|_{H^1(\Omega)} \leq \|u\|_{L^2(\Omega)} + 2h^\theta \|u\|_{H^{1+\theta}(\Omega)}$ . Combining these estimates, we obtain

$$\inf_{v \in V_N} \|u - v\|_{L^2(\Omega)} + h \|\nabla(u - v)\|_{L^2(\Omega)} \leq \left( h^{1+\theta} + (kh)^2 h^\theta \right) \|u\|_{H^{1+\theta}(\Omega)} + (kh)^2 \|u\|_{L^2(\Omega)},$$

which concludes the proof.  $\square$

**Lemma A.6.** *Let  $\beta \in [0, 1)$ . Then, for every  $p \in \mathbb{N}$  there exists a linear operator  $\pi_p : H_\beta^{2,2}(\hat{K}) \rightarrow \mathcal{P}_p$  that admits an “element-by-element construction” in the sense of [61, Def. 5.3] with the following approximation property:*

$$p \|u - \pi_p u\|_{L^2(\hat{K})} + \|u - \pi_p u\|_{H^1(\hat{K})} \leq C p^{-(1-\beta)} \|r^\beta \nabla^2 u\|_{L^2(\hat{K})},$$

where  $C > 0$  is independent of  $p$  and  $u$ .

*Proof.* Inspection of the proof of [61, Thm. B.4] shows that the operator  $\pi_p$  constructed there does in fact not depend on the regularity parameter  $s > 1$ . It has (as stated in [61, Thm. B.4]), the approximation property

$$p \|u - \pi_p u\|_{L^2(\hat{K})} + \|u - \pi_p u\|_{H^1(\hat{K})} \leq C p^{-(s-1)} \|u\|_{H^s(\hat{K})} \quad \forall u \in H^s(\hat{K}), \quad (\text{A.4})$$

if  $p \geq s - 1$ . Upon writing the Besov space  $B_{2,\infty}^s$  as an interpolation space  $B_{2,\infty}^s = (H^2(\hat{K}), H^1(\hat{K}))_{s-1,\infty}$  for  $s \in (1, 2)$ , we can infer for  $s \in (1, 2)$  from (A.4) the slightly stronger statement

$$p \|u - \pi_p u\|_{L^2(\hat{K})} + \|u - \pi_p u\|_{H^1(\hat{K})} \leq C p^{-(s-1)} \|u\|_{B_{2,\infty}^s(\hat{K})} \quad \forall u \in B_{2,\infty}^s(\hat{K}). \quad (\text{A.5})$$

Appealing to Lemma A.1 then yields

$$p \|u - \pi_p u\|_{L^2(\hat{K})} + \|u - \pi_p u\|_{H^1(\hat{K})} \leq C p^{-(s-1)} \|u\|_{H_\beta^{2,2}(\hat{K})}. \quad (\text{A.6})$$

We replace the full  $H_\beta^{2,2}(\hat{K})$  norm by the seminorm in the standard way by a compactness argument. Since  $H_\beta^{2,2}(\hat{K})$  is compactly embedded in  $H^1(\hat{K})$  (see, e.g., [77, Lemma 4.19]) one obtains  $\inf_{v \in \mathcal{P}_1} \|u - v\|_{H_\beta^{2,2}(\hat{K})} \leq C \|r^\beta \nabla^2 u\|_{L^2(\hat{K})}$ . The proof is completed by noting that (A.4) implies that  $\pi_p$  reproduces linear polynomials.  $\square$

## References

1. R. A. Adams. *Sobolev Spaces*. Academic Press, 1975.
2. M. Ainsworth, P. Monk, and W. Muniz. Dispersive and dissipative properties of discontinuous Galerkin finite element methods for the second-order wave equation. *J. Sci. Comput.*, 27(1-3):5–40, 2006.
3. Mark Ainsworth. Discrete dispersion relation for  $hp$ -version finite element approximation at high wave number. *SIAM J. Numer. Anal.*, 42(2):553–575, 2004.
4. Mark Ainsworth and Hafiz Abdul Wajid. Dispersive and dissipative behavior of the spectral element method. *SIAM J. Numer. Anal.*, 47(5):3910–3937, 2009.
5. R. J. Astley and P. Gamallo. Special short wave elements for flow acoustics. *Comput. Methods Appl. Mech. Engrg.*, 194(2-5):341–353, 2005.
6. A. K. Aziz, R. B. Kellogg, and A. B. Stephens. A two point boundary value problem with a rapidly oscillating solution. *Numer. Math.*, 53(1-2):107–121, 1988.
7. V. M. Babič and V. S. Buldyrev. *Short-wavelength diffraction theory*, volume 4 of *Springer Series on Wave Phenomena*. Springer-Verlag, Berlin, 1991. Asymptotic methods, Translated from the 1972 Russian original by E. F. Kuester.
8. I. Babuška and B.Q. Guo. The  $h - p$  version of the finite element method. Part 1: The basic approximation results. *Computational Mechanics*, 1:21–41, 1986.
9. I. Babuška, F. Ihlenburg, E. Paik, and S. Sauter. A generalized finite element method for solving the Helmholtz equation in two dimensions with minimal pollution. *Comput. Meth. Appl. Mech. Engrg.*, 128:325–360, 1995.
10. I. Babuška, R.B. Kellogg, and J. Pitkäranta. Direct and inverse error estimates for finite elements with mesh refinements. *Numer. Math.*, 33:447–471, 1979.
11. I. Babuška and J.E. Osborn. Eigenvalue problems. In *Handbook of Numerical Analysis, Vol. II*, pages 641–789. North Holland, 1991.
12. I. Babuška and S. Sauter. Is the pollution effect of the FEM avoidable for the Helmholtz equation considering high wave numbers? *SIAM Review*, 42:451–484, 2000.
13. L. Banjai and S. Sauter. A refined Galerkin error and stability analysis for highly indefinite variational problems. *SIAM J. Numer. Anal.*, 45(1):37–53, 2007.
14. A. H. Barnett and T. Betcke. An exponentially convergent nonpolynomial finite element method for time-harmonic scattering from polygons. *SIAM J. Sci. Stat. Comp.*, 32:1417–1441, 2010.
15. A. Bayliss, C.I. Goldstein, and E. Turkel. On accuracy conditions for the numerical computation of waves. *J. Comput. Physics*, 59:396–404, 1985.
16. T. Betcke, S. Chandler-Wilde, I. Graham, S. Langdon, and M. Lindner. Condition number estimates for combined potential integral operators in acoustics and their boundary element discretization. *Numer. Meths. PDEs*, 27:31–69, 2011.
17. James H. Bramble and Joseph E. Pasciak. Analysis of a finite element PML approximation for the three dimensional time-harmonic Maxwell problem. *Math. Comp.*, 77(261):1–10, 2008.
18. S.C. Brenner and L.R. Scott. *The mathematical theory of finite element methods*. Springer Verlag, 1994.
19. A. Buffa and P. Monk. Error estimates for the ultra weak variational formulation of the Helmholtz equation. *M2AN (Mathematical Modelling and Numerical Analysis)*, 42(6):925–940, 2008.
20. O. Cessenat and B. Després. application of the ultra-weak variational formulation to the 2d Helmholtz problem. *SIAM J. Numer. Anal.*, 35:255–299, 1998.
21. O. Cessenat and B. Després. Using plane waves as base functions for solving time harmonic equations with the ultra weak variational formulation. *J. Computational Acoustics*, 11:227–238, 2003.
22. S. Chandler-Wilde and I. Graham. Boundary integral methods in high frequency scattering. In B. Engquist, A. Fokas, E. Hairer, and A. Iserles, editors, *highly oscillatory problems*. Cambridge University Press, 2009.

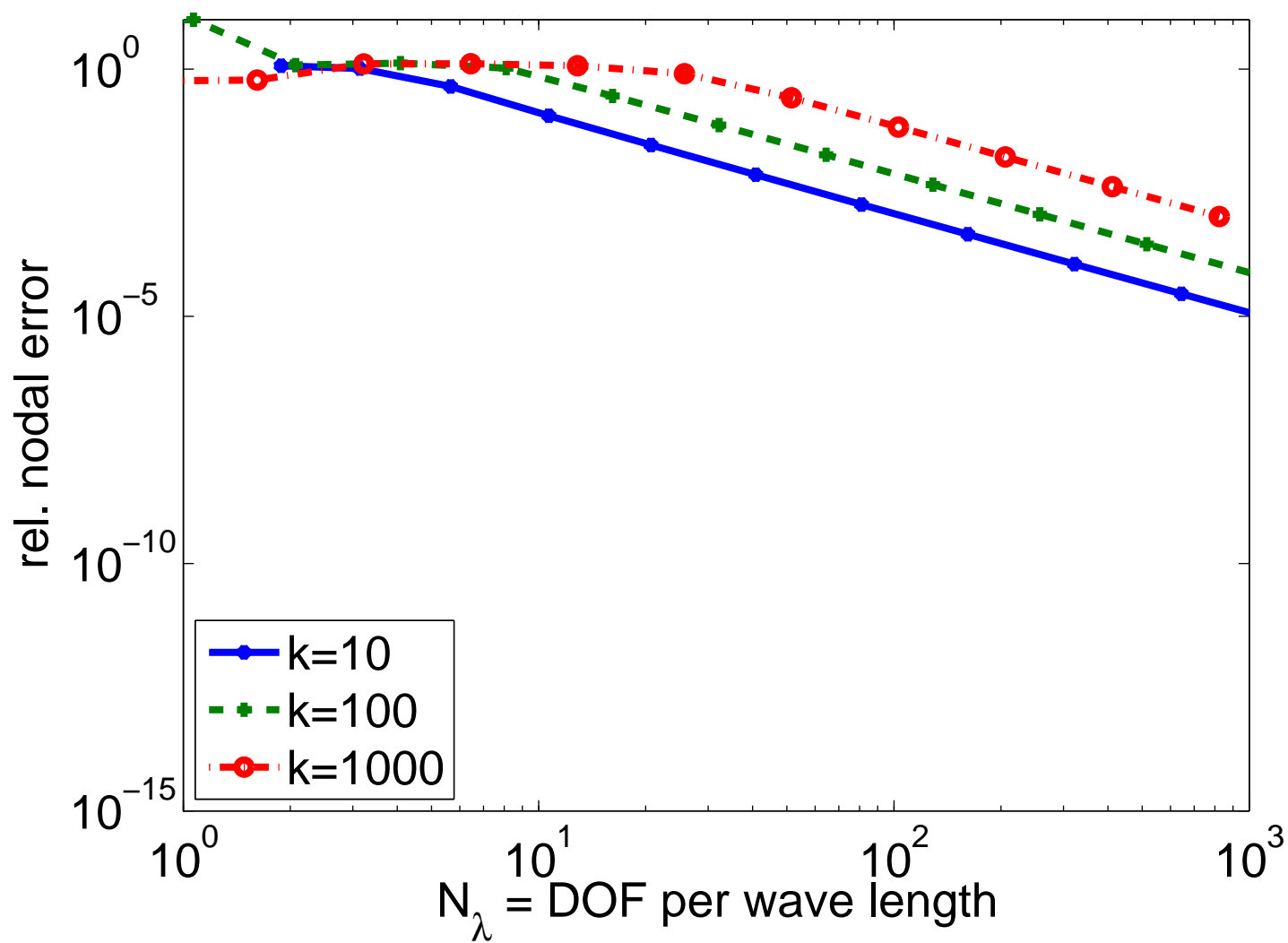
23. S.N. Chandler-Wilde and P. Monk. Wave-number-explicit bounds in time-harmonic scattering. *SIAM J. Math. Anal.*, 39:1428–1455, 2008.
24. Francis Collino and Peter Monk. The perfectly matched layer in curvilinear coordinates. *SIAM J. Sci. Comput.*, 19(6):2061–2090, 1998.
25. Peter Cummings and Xiaobing Feng. Sharp regularity coefficient estimates for complex-valued acoustic and elastic Helmholtz equations. *Math. Models Methods Appl. Sci.*, 16(1):139–160, 2006.
26. L. Demkowicz and K. Gerdes. Convergence of the infinite element methods for the Helmholtz equation in separable domains. *Numer. Math.*, 79(1):11–42, 1998.
27. A. Deraemaeker, I. Babuška, and P. Bouillard. Dispersion and pollution of the FEM solution for the Helmholtz equation in one, two and three dimensions. *Int. J. Numer. Meth. Eng.*, 46(4):471–499, 1999.
28. B. Engquist and L. Ying. Sweeping preconditioner for the Helmholtz equation: Hierarchical matrix representation. Technical report, Dept. of Mathematics, UT Austin, 2010.
29. Björn Engquist and Olof Runborg. Computational high frequency wave propagation. *Acta Numer.*, 12:181–266, 2003.
30. Yogi A. Erlangga. Advances in iterative methods and preconditioners for the Helmholtz equation. *Arch. Comput. Methods Eng.*, 15(1):37–66, 2008.
31. C. Farhat, I. Harari, and U. Hetmaniuk. A discontinuous Galerkin method with Lagrange multipliers for the solution of Helmholtz problems in the mid-frequency regime. *Comp. Meth. Appl. Mech. Eng.*, 192:1389–1419, 2003.
32. Charbel Farhat, Radek Tezaur, and Paul Weidemann-Goiran. Higher-order extensions of a discontinuous Galerkin method for mid-frequency Helmholtz problems. *Internat. J. Numer. Methods Engrg.*, 61(11):1938–1956, 2004.
33. Xiaobing Feng and Haijun Wu. Discontinuous Galerkin methods for the Helmholtz equation with large wave number. *SIAM J. Numer. Anal.*, 47(4):2872–2896, 2009.
34. Xiaobing Feng and Haijun Wu. *hp*-discontinuous Galerkin methods for the Helmholtz equation with large wave number. *Math. Comp.*, 80:1997–2024, 2011.
35. Xiaobing Feng and Yulong Xing. Absolutely stable local discontinuous Galerkin methods for the Helmholtz equation with large wave number. Technical report, 2010. [arXiv:1010.4563v1](https://arxiv.org/abs/1010.4563v1) [math.NA].
36. C. Gittelsohn, R. Hiptmair, and I. Perugia. Plane wave discontinuous Galerkin methods. *M2AN (Mathematical Modelling and Numerical Analysis)*, 43:297–331, 2009.
37. Dan Givoli. *Numerical methods for problems in infinite domains*, volume 33 of *Studies in Applied Mechanics*. Elsevier Scientific Publishing Co., Amsterdam, 1992.
38. P. Grisvard. *Elliptic Problems in Nonsmooth Domains*. Pitman, 1985.
39. Isaac Harari. A survey of finite element methods for time-harmonic acoustics. *Comput. Methods Appl. Mech. Engrg.*, 195(13-16):1594–1607, 2006.
40. Isaac Harari and Thomas J. R. Hughes. Galerkin/least-squares finite element methods for the reduced wave equation with nonreflecting boundary conditions in unbounded domains. *Comput. Methods Appl. Mech. Engrg.*, 98(3):411–454, 1992.
41. U. Hetmaniuk. Stability estimates for a class of Helmholtz problems. *Commun. Math. Sci.*, 5(3):665–678, 2007.
42. R. Hiptmair, A. Moiola, and I. Perugia. Approximation by plane waves. Technical Report 2009-27, Seminar für Angewandte Mathematik, ETH Zürich, 2009.
43. R. Hiptmair, A. Moiola, and I. Perugia. Plane wave discontinuous Galerkin methods for the 2d Helmholtz equation: analysis of the *p*-version. *SIAM J. Numer. Anal.*, 49:264–284, 2011.
44. Thorsten Hohage and Lothar Nannen. Hardy space infinite elements for scattering and resonance problems. *SIAM J. Numer. Anal.*, 47(2):972–996, 2009.
45. T. Huttunen, P. Gamallo, and R. J. Astley. Comparison of two wave element methods for the Helmholtz problem. *Comm. Numer. Methods Engrg.*, 25(1):35–52, 2009.
46. T. Huttunen and P. Monk. The use of plane waves to approximate wave propagation in anisotropic media. *J. Computational Mathematics*, 25:350–367, 2007.
47. F. Ihlenburg. *Finite Element Analysis of Acoustic Scattering*, volume 132 of *Applied Mathematical Sciences*. Springer Verlag, 1998.



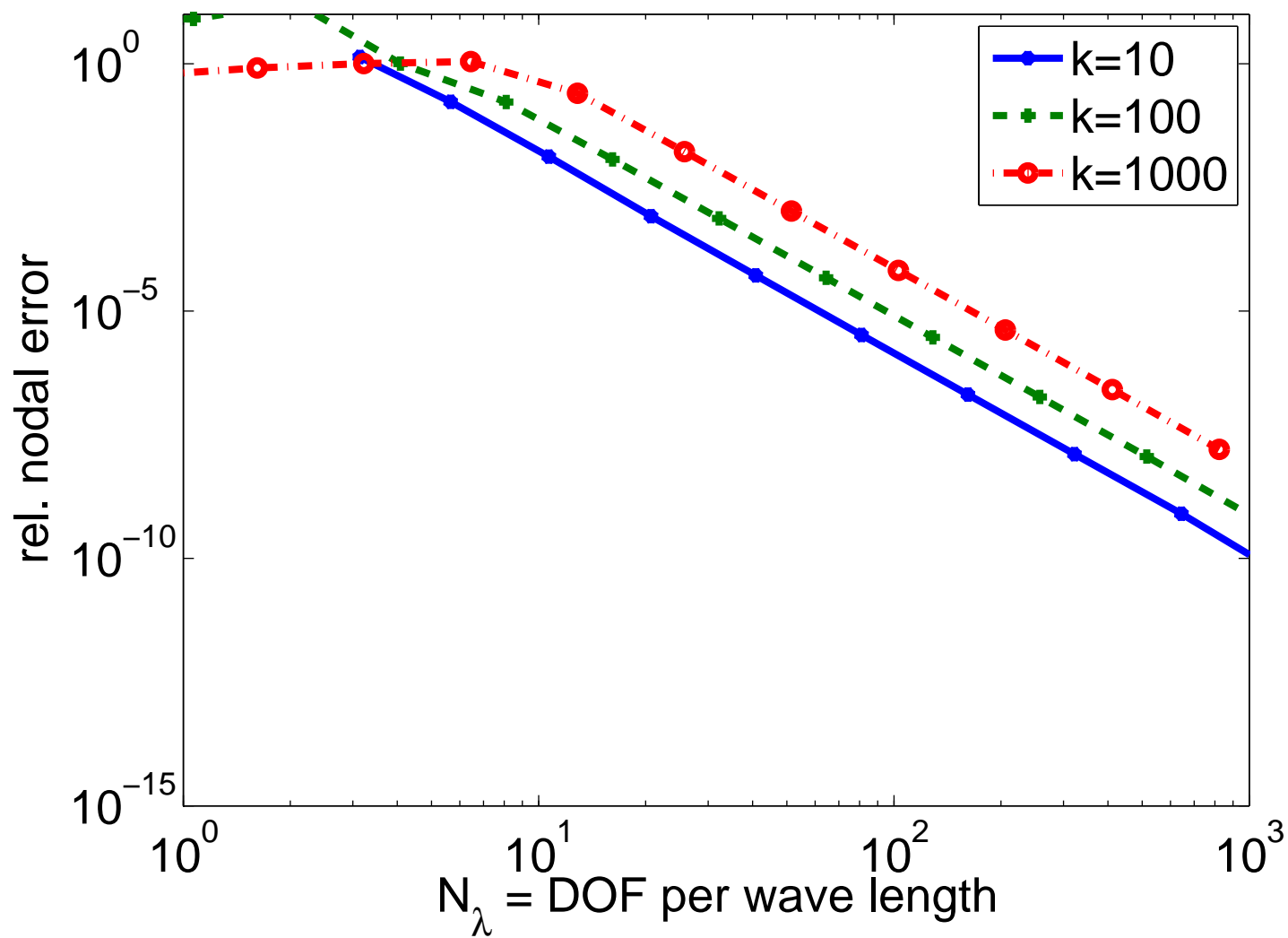
48. F. Ihlenburg. Sound in vibrating cabins: Physical effects, mathematical description, computational simulation with FEM. In G. Sandberg and P. Ohayon, editors, *Computational Aspects of Structural Acoustics and Vibration*, pages 102–170. Springer Verlag, 2009.
49. F. Ihlenburg and I. Babuška. Finite element solution to the Helmholtz equation with high wave number. Part II: The *hp*-version of the fem. *SIAM J. Numer. Anal.*, 34:315–358, 1997.
50. O. Laghrouche and P. Bettess. Solving short wave problems using special finite elements; towards an adaptive approach. In J. Whiteman, editor, *Mathematics of Finite Elements and Applications X*, pages 181–195. Elsevier, 2000.
51. O. Laghrouche, P. Bettess, and J. Astley. Modelling of short wave diffraction problems using approximation systems of plane waves. *Internat. J. Numer. Meths. Engrg.*, 54:1501–1533, 2002.
52. R. Leis. *Initial Boundary Value Problems in Mathematical Physics*. Teubner, Wiley, 1986.
53. Z. C. Li. The Trefftz method for the Helmholtz equation with degeneracy. *Appl. Numer. Math.*, 58(2):131–159, 2008.
54. M. Löhndorf and J.M. Melenk. Wavenumber-explicit *hp*-BEM for high frequency scattering. Technical Report 02/2010, Institute for Analysis and Scientific Computing, TU Wien, 2010.
55. Teemu Luostari, Tomi Huttunen, and Peter Monk. Plane wave methods for approximating the time harmonic wave equation. In *Highly oscillatory problems*, volume 366 of *London Math. Soc. Lecture Note Ser.*, pages 127–153. Cambridge Univ. Press, Cambridge, 2009.
56. J. M. Melenk. *On Generalized Finite Element Methods*. PhD thesis, University of Maryland, 1995.
57. J.M. Melenk. *hp finite element methods for singular perturbations*, volume 1796 of *Lecture Notes in Mathematics*. Springer Verlag, 2002.
58. J.M. Melenk. On approximation in meshless methods. In J. Blowey and A. Craig, editors, *Frontier in Numerical Analysis, Durham 2004*, pages 65–141. Springer Verlag, 2005.
59. J.M. Melenk. Mapping properties of combined field Helmholtz boundary integral operators. Technical Report 01/2010, Institute for Analysis and Scientific Computing, TU Wien, 2010.
60. J.M. Melenk and I. Babuška. The partition of unity finite element method: Basic theory and applications. *Comput. Meth. Appl. Mech. Engrg.*, 139:289–314, 1996.
61. J.M. Melenk and S. Sauter. Convergence analysis for finite element discretizations of the Helmholtz equation with Dirichlet-to-Neumann boundary conditions. *Math. Comp.*, 79:1871–1914, 2010.
62. J.M. Melenk and S. Sauter. Wavenumber explicit convergence analysis for finite element discretizations of the Helmholtz equation. *SIAM J. Numer. Anal.*, 49:1210–1243, 2011.
63. A. Moiola. Approximation properties of plane wave spaces and application to the analysis of the plane wave discontinuous Galerkin method. Technical Report 2009-06, Seminar für Angewandte Mathematik, ETH Zürich, 2009.
64. P. Monk, J. Schöberl, and A. Sinwel. Hybridizing Raviart-Thomas elements for the Helmholtz equation. *Electromagnetics*, 30:149–176, 2010.
65. P. Monk and D.Q. Wang. A least squares methods for the Helmholtz equation. *Comput. Meth. Appl. Mech. Engrg.*, 175:121–136, 1999.
66. C. S. Morawetz and D. Ludwig. An inequality for the reduced wave operator and the justification of geometrical optics. *Comm. Pure Appl. Math.*, 21:187–203, 1968.
67. J. Nečas. *Les méthodes directes en théorie des équations elliptiques*. Masson, 1967.
68. Pablo Ortiz. Finite elements using a plane-wave basis for scattering of surface water waves. *Philos. Trans. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci.*, 362(1816):525–540, 2004.
69. E. Perrey-Debain, O. Laghrouche, P. Bettess, and J. Trevelyan. Plane-wave basis finite elements and boundary elements for three-dimensional wave scattering. *Philos. Trans. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci.*, 362(1816):561–577, 2004.
70. B. Pluymers, B. van Hal, D. Vandepitte, and W. Desmet. Trefftz-based methods for time-harmonic acoustics. *Arch. Comput. Methods Eng.*, 14(4):343–381, 2007.
71. S. I. Pohožaev. On the eigenfunctions of the equation  $\Delta u + \lambda f(u) = 0$ . *Dokl. Akad. Nauk SSSR*, 165:36–39, 1965.
72. S.A. Sauter. A Refined Finite Element Convergence Theory for Highly Indefinite Helmholtz Problems. *Computing*, 78(2):101–115, 2006.

73. S.A. Sauter and C. Schwab. *Boundary element methods*. Springer Verlag, 2010.
74. Alfred H. Schatz. An observation concerning Ritz-Galerkin methods with indefinite bilinear forms. *Math. Comp.*, 28:959–962, 1974.
75. J. Schöberl. *Finite Element Software Netgen/NGSolve version 4.13*. <http://sourceforge.net/projects/ngsolve/>.
76. J. Schöberl. NETGEN - an advancing front 2d/3d-mesh generator based on abstract rules. *Computing and Visualization in Science*, 1(1):41–52.
77. C. Schwab. *p- and hp-Finite Element Methods*. Oxford University Press, 1998.
78. E.A. Spence, S.N. Chandler-Wilde, I.G. Graham, and V.P. Smyshlyaev. A new frequency-uniform coercive boundary integral equation for acoustic scattering. *Comm. Pure Appl. Math.*, 60(10):1384–1415, 2011.
79. E.M. Stein. *Singular integrals and differentiability properties of functions*. Princeton University Press, 1970.
80. M. Stojek. Least squares Trefftz-type elements for the Helmholtz equation. *Internat. J. Numer. Meths. Engrg.*, 41:831–849, 1998.
81. Theofanis Strouboulis, Ivo Babuška, and Realino Hidajat. The generalized finite element method for Helmholtz equation: theory, computation, and open problems. *Comput. Methods Appl. Mech. Engrg.*, 195(37–40):4711–4731, 2006.
82. R. Tezaur and C. Farhat. Three-dimensional discontinuous Galerkin elements with plane waves and Lagrange multipliers for the solution of mid-frequency Helmholtz problems. *Internat. J. Numer. Meths. Engrg.*, 66:796–815, 2006.
83. L.L. Thompson. A review of finite element methods for time-harmonic acoustics. *Journal Acoustical Society America*, 119:1315–1330, 2006.

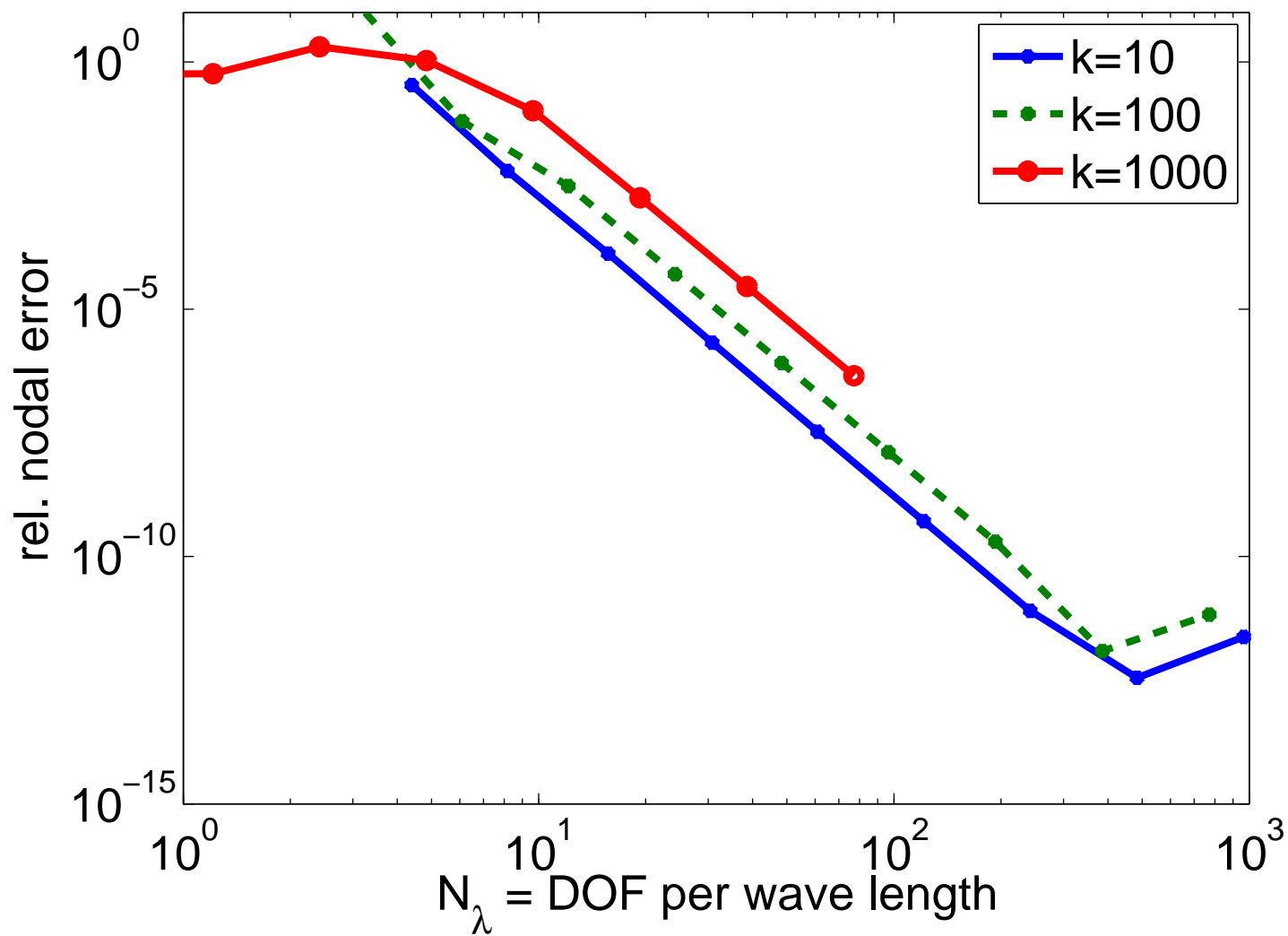
$$p=1; -u'' - k^2 u = 1$$



$$p=2; -u'' - k^2 u = 1$$



$$p=3; -u'' - k^2 u = 1$$



$$p=4; -u'' - k^2 u = 1$$

